

# INIA Uruguay 33

Taller de Big Data de arroz

Everett Wetchler



BAYES IMPACT



BAYES IMPACT

Misión: avanzar causas sociales con datos, algoritmos, y software.

- ONG, fundado en 2014 en San Francisco, California
- Ingenieros de sistema, estadísticos, y científicos de datos.
- Trabajamos con otros ONGs, gobiernos, y científicos, como consultorio o colegas.
- Nos enfocamos en problemas con gran impacto para seres humano.

# Objetivos

- Mostrar metodos y herramientas mas de hacer modelos de machine learning.
- Python vs R
- Metodos de limpieza y exploracion de datos
- Correlaciones basicos
- Hacer un poco de "feature engineering" -- la ingeniería de variables

# iPython Notebook

- Editor interactivo de browser, originalmente para ciencia de datos de python
- Incluye codigos y graficos en linea
- Naturalmente documenta su trabajo, como tranformo los datos, cuales codigos hicieron los graficos que estamos viendo, etc.
- Organizacion del proyecto, documenta pasos.
- **Se puede usar con R!**

# jupyter aprender\_r (unsaved changes)

File Edit View Insert Cell Kernel Help

Cell Toolbar: None

```
In [5]: 2 + 4
```

```
Out[5]: 6
```

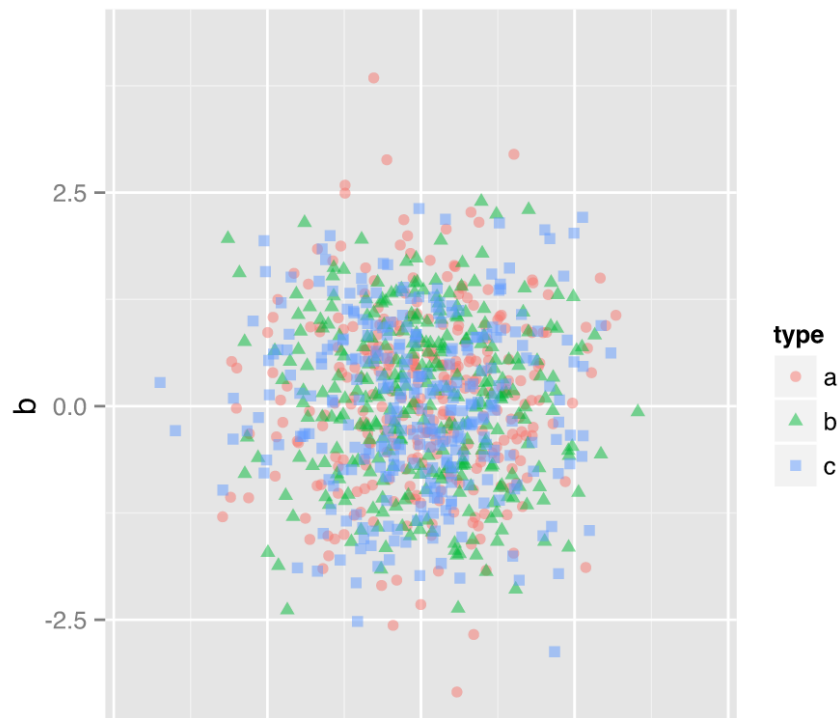
```
In [1]: library(ggplot2)
```

```
In [4]: points = 300
types = 3
N = points * types
df <- data.frame(type=rep(letters[1:types], each = points),
                  a=rnorm(N), b=rnorm(N))
```

```
In [3]: xabsmax = max(sapply(df$a, abs))*1.1
yabsmax = max(sapply(df$b, abs))*1.1
ggplot(df) +
  aes(x = a, y = b,
      xmin=-xabsmax, xmax=xabsmax,
      ymin=-yabsmax, ymax=yabsmax,
      shape=type,
      color=type) +
  geom_point(size=2, alpha=0.5)
```

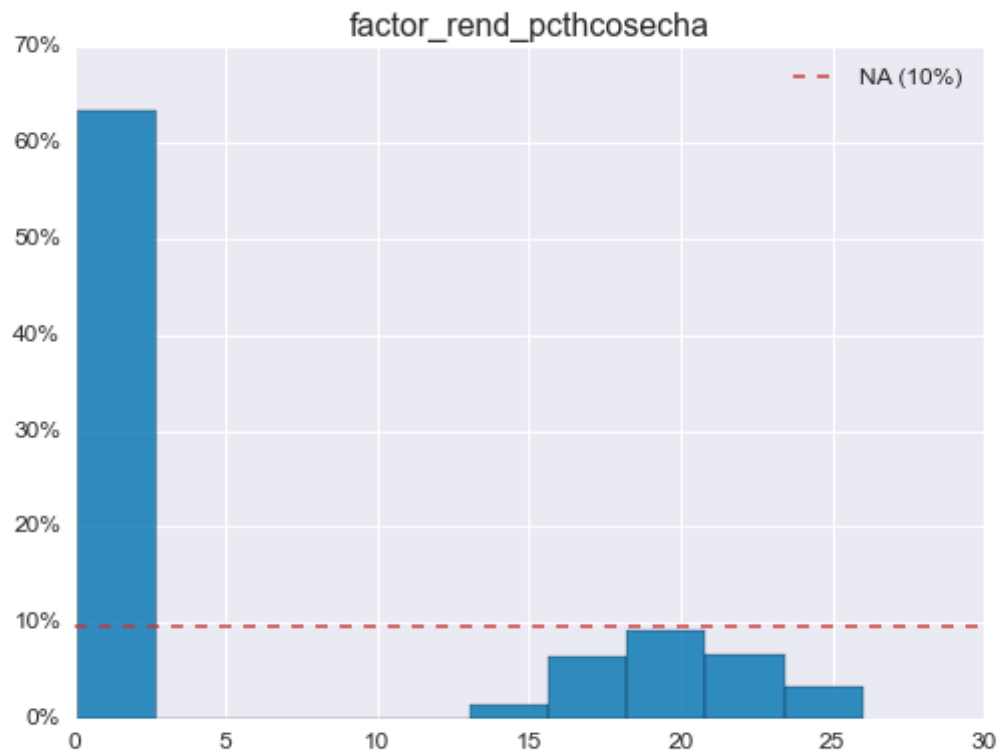
```
Error in if (args[[1]]$name == "C title" && !is.null(args[[2]])) {: missing value where TRUE/FALSE needed
```

```
In [7]: ggplot(df) +  
  aes(x = a, y = b,  
      xmin=-xabsmax, xmax=xabsmax, ymin=-yabsmax, ymax=yabsmax, shape=type, color=type) +  
  geom_point(size=2, alpha=0.5)  
  
Error in if (args[[1]]$name == "C_title" && !is.null(args[[2]])) {: missing value where TRUE/FALSE needed
```



# Limpieza de datos

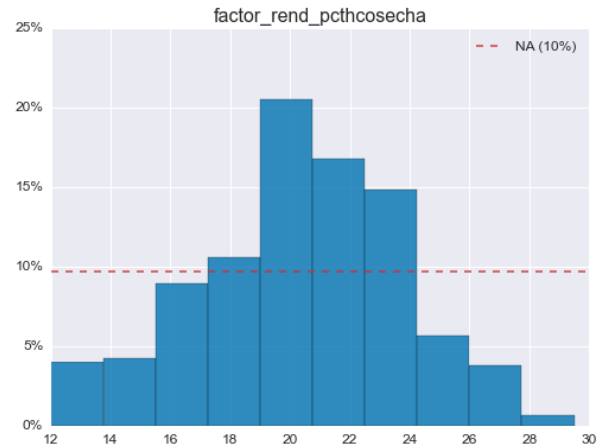
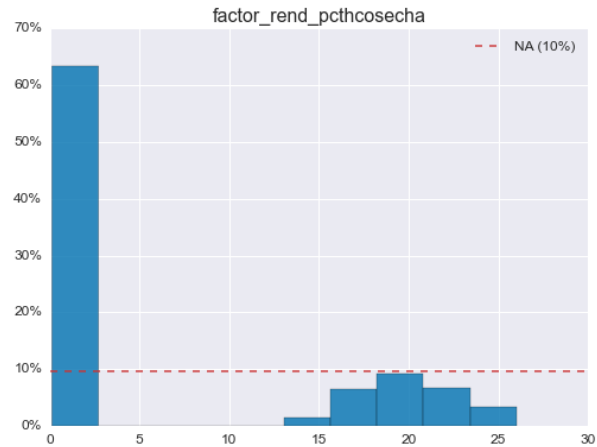
- No cambio los datos originales!
- Ver visualmente los datos en Excel... pero quiero data frame
- Codigos para la limpieza
  - input: datos originales (csv, excel, ...)
  - output: csv limpio
- Codigos que automaticamente hacen histograms y bars para explorar todos los columnos visualmente
- Proceso iterativo: Encontrar problema con datos -> crear codigos para limpiarlos -> averiguar que mejore el problema -> encontrar otro problema....
- Ejemplo vivo





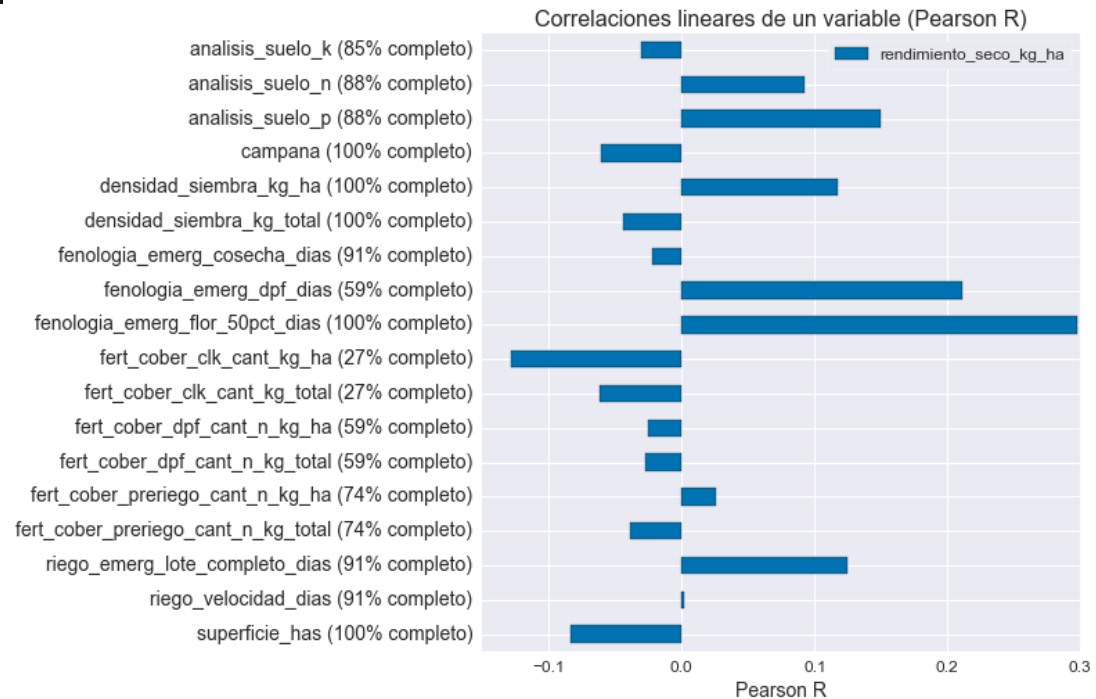
	AM	AN	AO	
	Fenologia_Emerg-Cosecha_Dias	Factor_Rend_%Hcosecha	Rendimiento_seco_kg/ha	
	0	41	0	
	136	21%	7877.25	
	131	23%	8597.19	
	131	18%	8503.07	
	138	19%	7815.19	
	139	19%	8945.95	
	132	18%	8503.07	
	148	14%	4953.83	
	148	14%	5073.54	
	147	14%	4834.69	
	120	20%	7063.19	
	147	14%	4444.51	
	127	21%	6236.61	
	127	21%	6241.53	
	143	1650.0%	7524.60	
	143		7524.60	
	143	1500.0%	8178.81	
	143		8178.81	
	142	1550.0%	7205.30	
	141	1550.0%	7205.30	
	141	1500.0%	7314.24	
	121	2050.0%	8386.82	
	121		8386.82	
	120	2400.0%	7973.65	
	131	2580.0%	10905.34	

# Limpieza

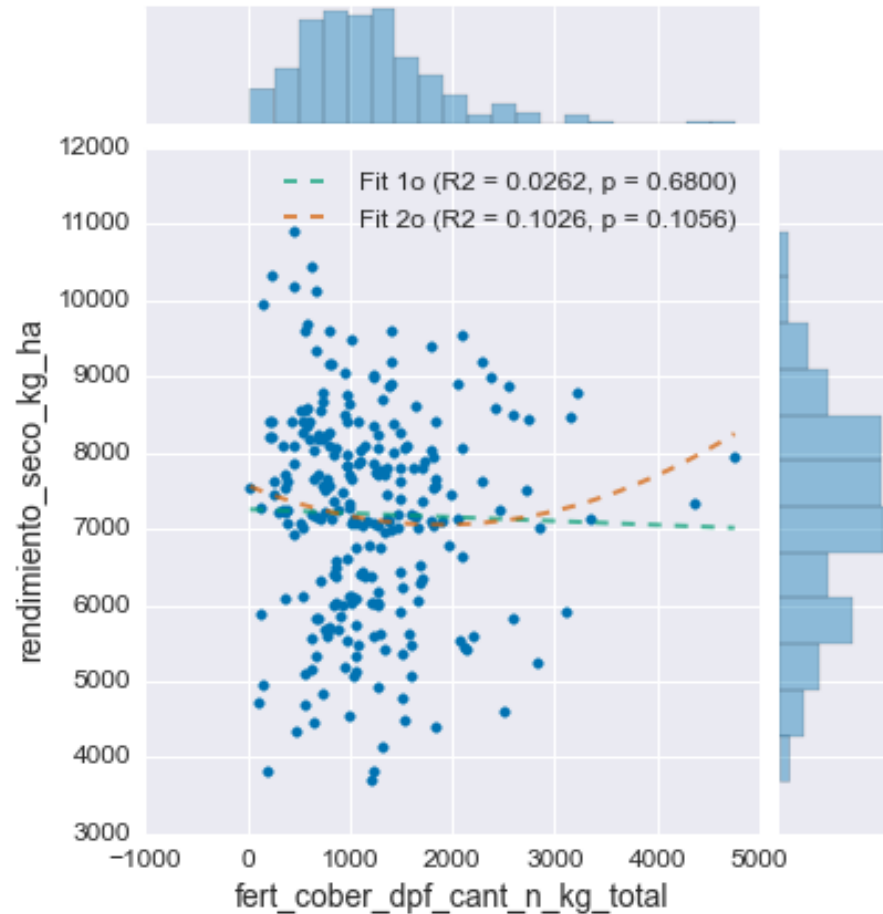


# Correlaciones basicos

- Ya estan bastante limpios los datos
- `DataFrame.corr()` -- todos los correlaciones lineales entre variables (Pearson o qualquiera)
- Bar para visualizarlos
- Mira: siempre incluyo mediciones de datos vacios
- El paquete “seaborn” hace graficos guapos, similar a ggplot. Viene de estudiante PhD de Stanford.
- `seaborn.jointplot()`
- “numpy” (otro paquete) para hacer lineas y cuervas de fit

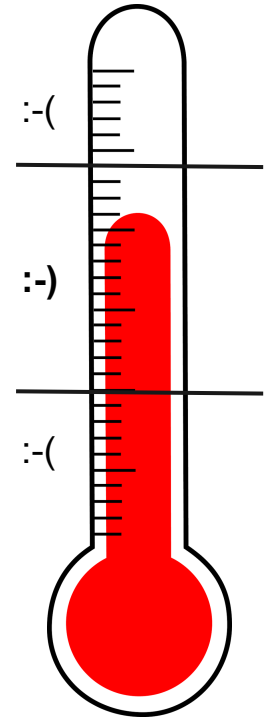


# Ejemplo



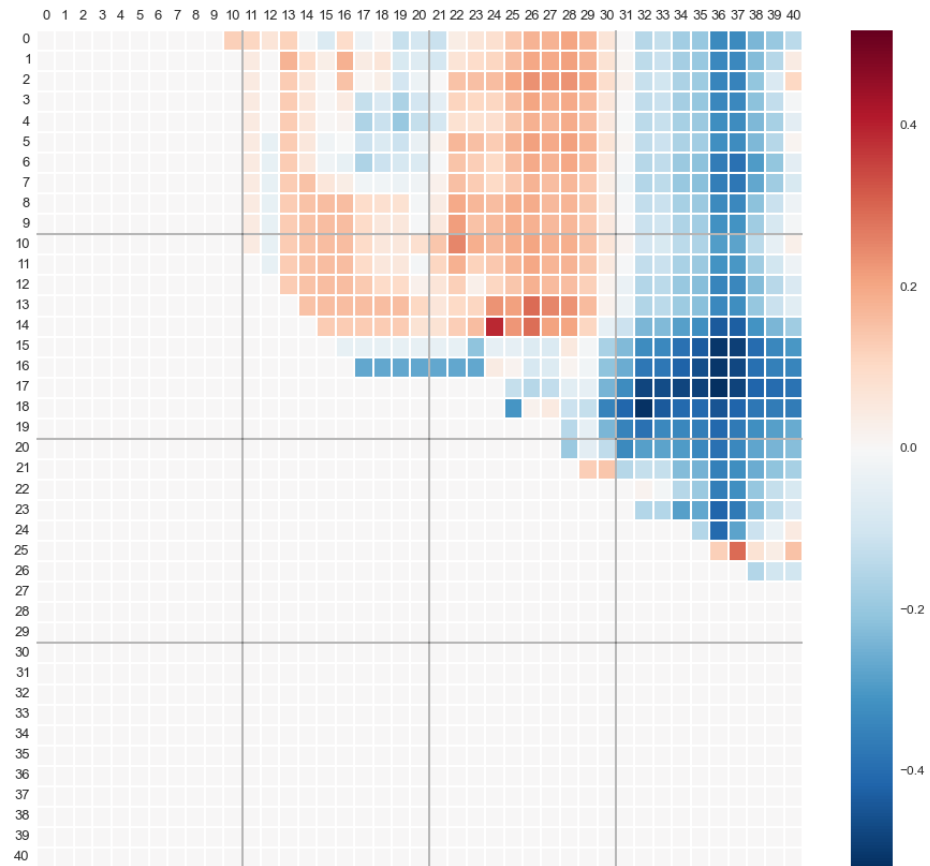
# Ingengeria de variables

- El mejor variable que ya tenemos explica unos ~30% de la variabilidad.
- No incluye datos de clima.
- Cuales variables climas debemos hacer? “Dias con  $T_{max} < 35$ ”?  $< 34$ ?  $T_{avg} > 20$ ?
- Teoria: hay un alcance critica que predice el exito del arroz. Mas dias durante el cutivo entre este alcance, mejor rienda.
  - 20-30? 20-32? 18-35?
  - Buscamos el mejor alcance

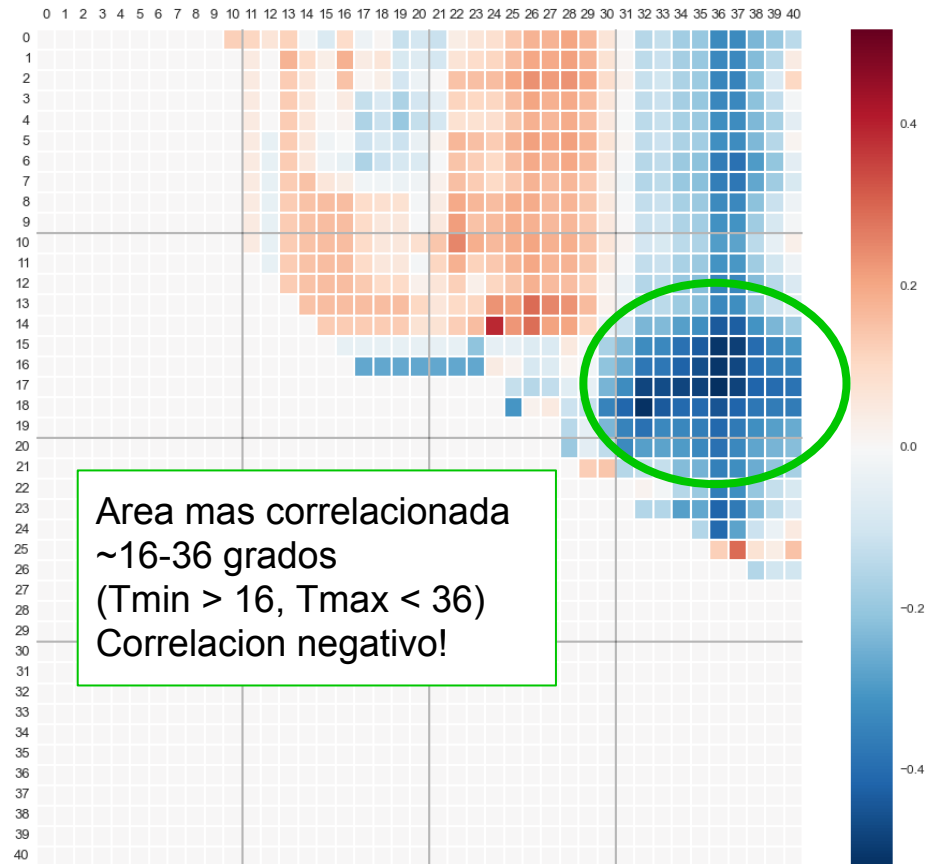


# Solucion

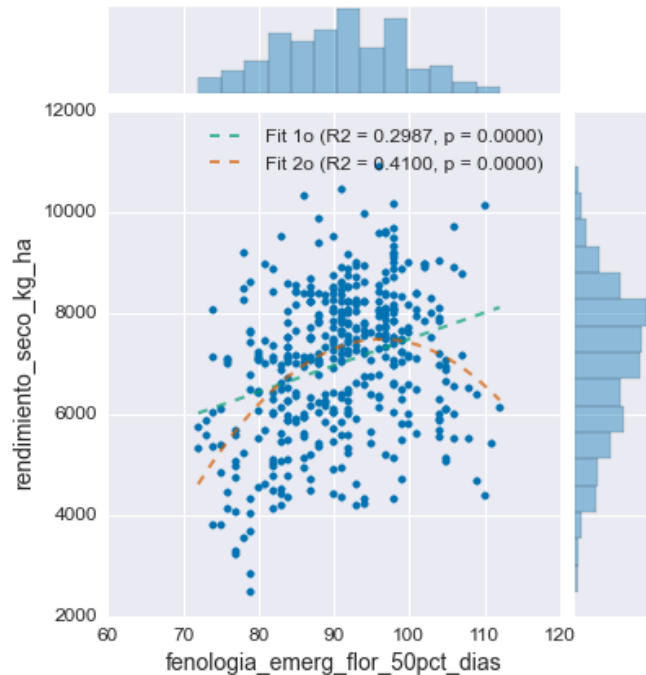
- Probar todas posibilidades ( $[0,1]$ ,  $[0,2]$ , ...,)
- Por cada alcance, calcula, para cada chacra, la porcentaje de los dias cultivos cuando la clima se mantenaba entre ellos todo el dia.
  - e.g.  $[20,30]$ , % dias con  $t_{max} < 30$  y  $t_{min} > 20$
- Medir cuanta variabilidad explicaria un variable asi ( $R^2$ )
- Resultado:







# Otra vista



**Mejor  
explicador**

**<- Antes**

**Despues ->**

Efecto de alcances críticas de temperatura (16°C-36°C)

