Other ML related activites

# Complex Social Science Gateway – a tool for cross-cultural analysis in R

Select dataset,
Select varialbes,
Submit analysis
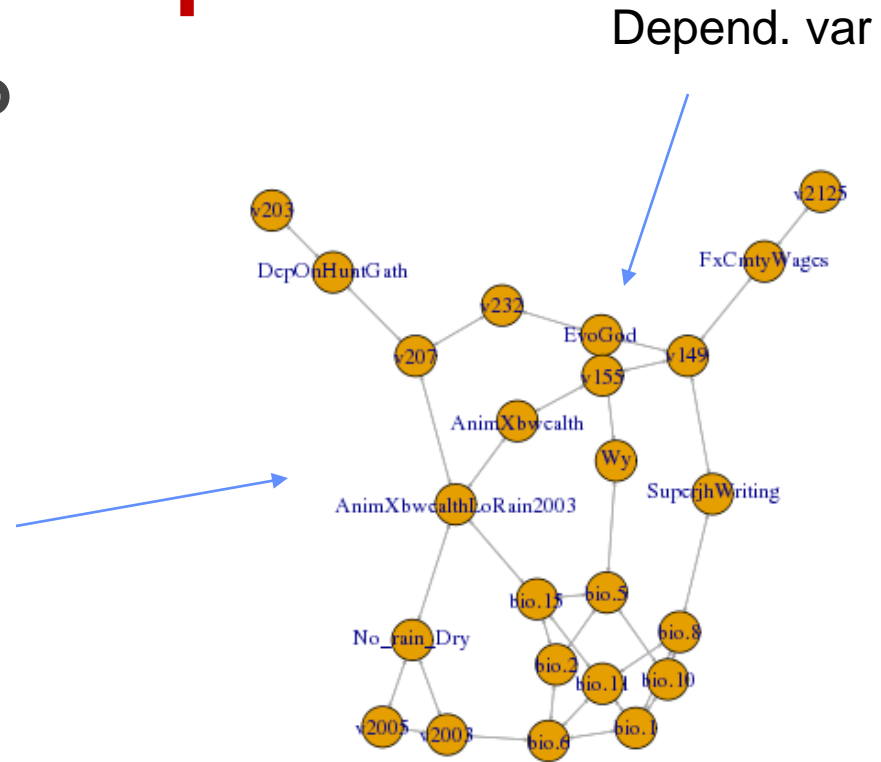
http://socscicom
pute.ss.uci.edu/

(but moving soon)

# R Analysis options

- **Two-stage least squares to handle spatially correlated errors (OLS, logit, multinomial logit)**

- **Bootstrap sampling of Bayesian network (package bnlearn) to confirm OLS effects, or suggest other moderating/mediating effects**

Depend. var

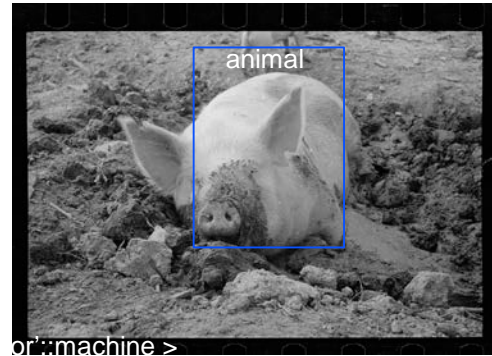# In a nutshell: Bayesian Network captures probabilistic dependencies between variables



AND:
Bootstrap samples determine most likely edges

# Image Analysis of Rural Photography 175K war and depression era photos extracting features for datamining

**Title:**
"Destitute pea pickers in California.
Mother of seven children."



**Histogram of Gradients**
CellSize = [32 32]
Feature length = 10260

For each pixel in a cell, take filters:
[-1 0 1]
And
[-1
  0
  1]
Take weighted average and bin into 9 orientations; the bin frequency is like magnitude

**Title:**

”Destitute pea pickers in California.
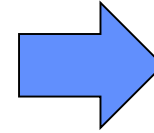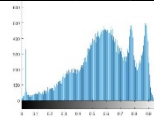Mother of seven children.”



CellSize = [32 32]
Feature length = 10260



Take all orientations, at different scales, as 1 big vector, and feed into classifier trained to recognize Face.

**Title:**

"Destitute pea pickers in California. Mother of seven children." By D. Lange, 1936, California, [metadata]

## Metadata processing:

- **Parse and tag speech (using Stanford NLP tools, word ontologies, in Python NLP toolkit)**
- **Several words identify 'person'**

- **SQL: give me all pictures by Lange with possible 'person' and num_faces > 0**

# Early 20th century, ~15k prison Bertillon id cards extracting information

Segment,
binarize,
denoise



extract field and cell



(word spotting) Get profile and
compare to known templates

# Linear to Logistic to Neural Network model

- $yi = bo * 1 + b_1 * xi_1 + b_2 * xi_2 \ldots = $ **B*X**

- **Squash** $bo * 1 + b_1 * xi_1$ **to 0,1 range using Logistic Function:**



B*X=0 then y=0.5

# Logistic Regression to Neural Networks

- **Use several squash functions (hidden layer)**



- **Take further combinations (output layer)**



- **More powerful but more complex**

many parameters, many options, needs more training

# organize connections into cells, add layers (deepen), add special pooling operations at some layers - you get a convolution network

# SciKit python package has a convolution neural network

```
nn2 = Classifier(
    layers=[
    Convolution("Rectifier", channels=numch, kernel_shape=(10,10),pool_shape=(2,2)),
    Convolution("Rectifier", channels=numch, kernel_shape=(6,6),pool_shape=(4,4)),
    Layer("Sigmoid",units=numalpha2do*4),
    Layer("Sigmoid",units=numalpha2do*2)
    ],
    verbose=False,
learning_rate=0.001,valid_set=(Xtrain,Ytrain),
n_iter=myiter)
nn2.fit(Xtrain,Ytrain)
```

# Topic Modelling with Latent Dirichlet Allocation

- **Each circle is 1 word occurrence**
- **2 topics (filled/empty circles), 15 docun**
- **Initially random assignments**



- **After learning, topics are well formed**

# LDA optimization

- **Start with initial guess of topic=$t$, and parameters**
- **Repeat:**

  Compute the expected value of word=$w$

  Compute the parameters that maximize likelihood L of *t given w*


  *Parameters are estimated from word/topic counts*


  *With each iteration, objective function L goes up*

# Topic Modelling with Latent Dirichlet Allocation on HPC

- **R LDA package: wraps C programs for Gibbs sampling or EM**

- **Mallet: Gibbs sampling      multicore, java code**

- **Spark LDA:  EM**

- **Asymptotic Distributed LDA :  MPI  based, no bells&whistles**

- **Example Case:**

articles from post WWII journals



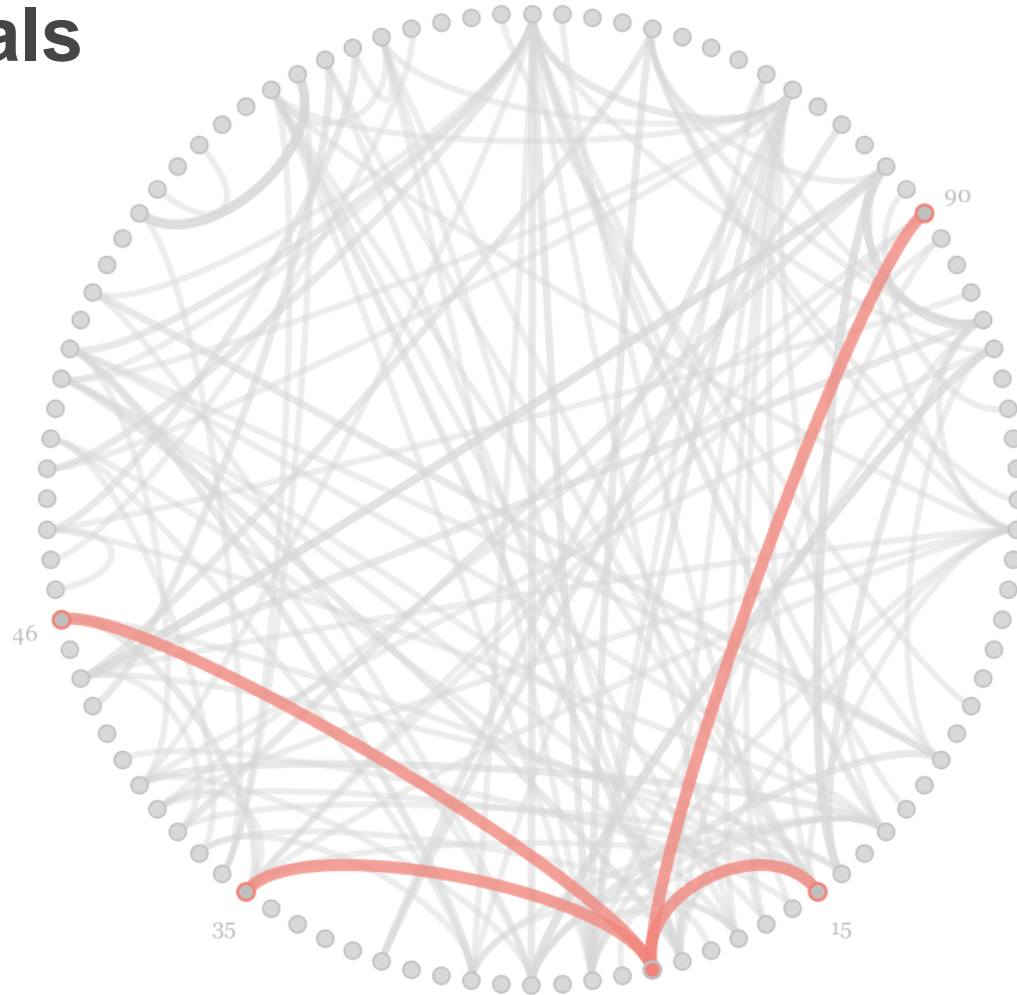| topic21 | topic35 | topic46 | topic90 | topic15 |
|---|---|---|---|---|
| worker | county | farm | house | work |
| labor | city | land | family | day |
| employment | state | farmer | area | labor |
| job | town | acre | home | time |
| percent | york | crop | city | pay |
| industry | public | agricultural | build | hour |
| defense | district | county | live | week |
| service | mayor | family | room | wage |
| work | relief | agriculture | unit | make |
| increase | local | cotton | community | month |
| employ | population | labor | project | condition |
| train | person | state | income | year |
| unemployment | citizen | area | rend | employ |
| wage | community | migrant | move | find |
| employee | place | rural | neighborhood | case |
| department | residence | year | facility | money |
| employer | resident | tenant | urban | service |
| occupation | board | migration | low | receive |
| rate | settlement | make | occupy | leave |
| number | large | large | resident | care |
| industrial | welfare | grower | lodge | employer |
| production | number | small | condition | order |
| earnings | part | california | apartment | good |
| woman | aid | camp | neighborhoods | require |

- **Example Case:**

**Sample topic plot (tree map)**