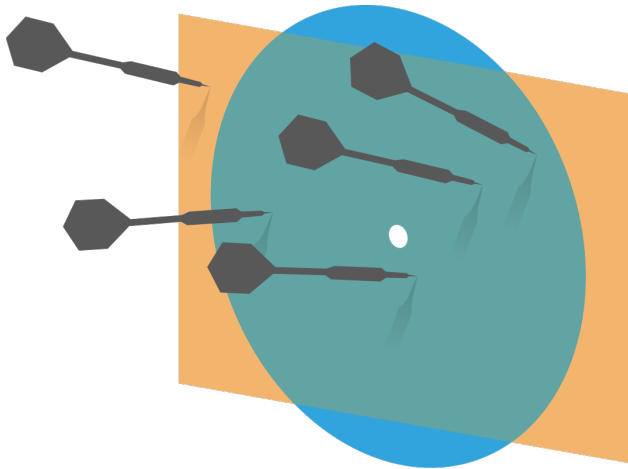


Probability and Statistics for Computer Science



“All models are wrong, but some models are useful” --- George Box

Credit: wikipedia

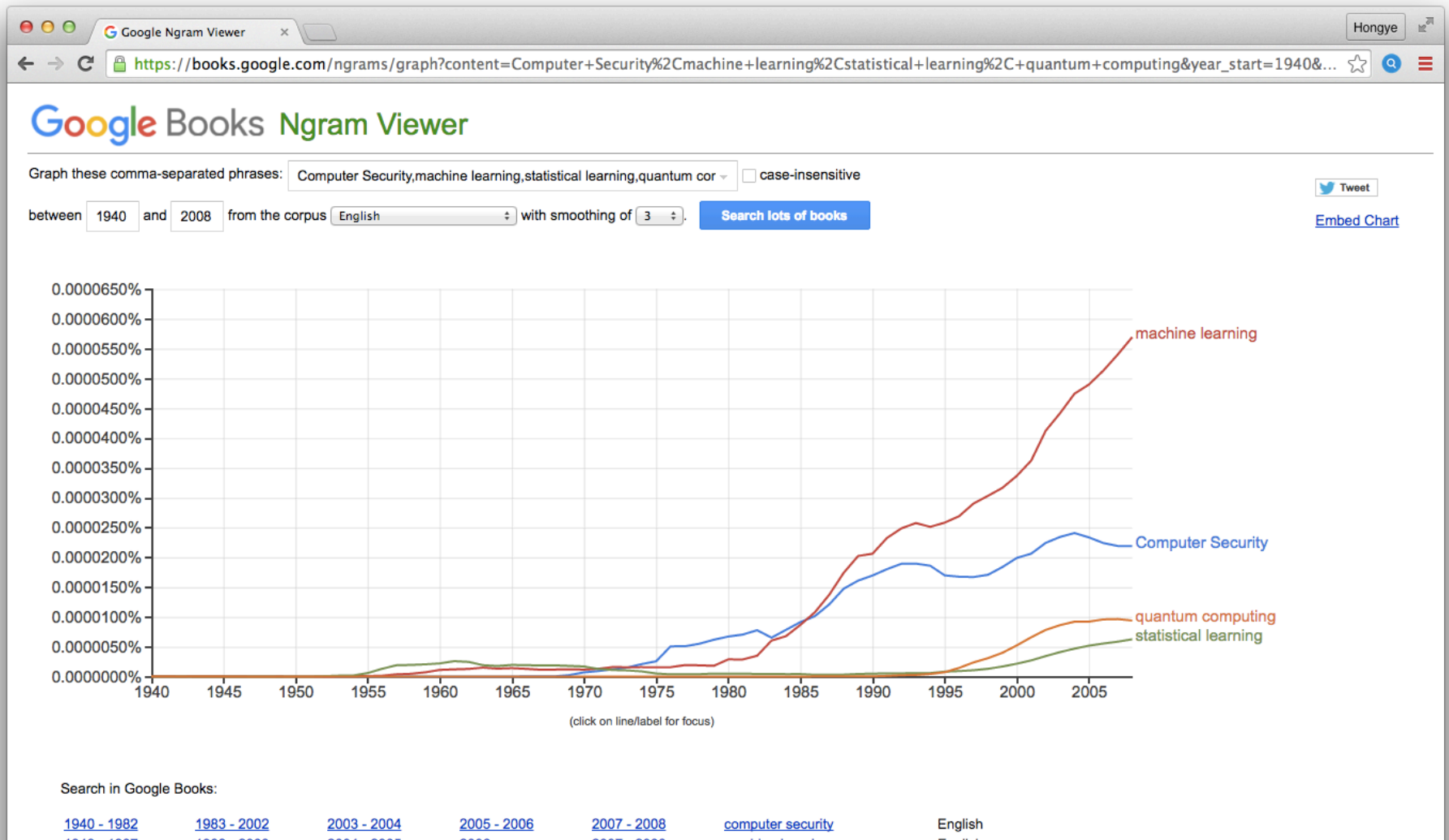
Last time

- ✱ Stochastic Gradient Descent
- ✱ Naïve Bayesian Classifier

Objectives

- * Linear regression
 - * The problem
 - * The least square solution
 - * The training and prediction
 - * The R-squared for the evaluation of the fit.

Some popular topics in Ngram



Regression models are Machine learning methods

- ✱ Regression models have been around for a while
- ✱ Dr. Kelvin Murphy's Machine Learning book has 3+ chapters on regression

Wait, have we seen the linear regression before?

It's about *Relationship* between data features

- ✱ Example: Is the height of people related to their weight?

IDNO	BODYFAT	DENSITY	AGE	WEIGHT	HEIGHT
1	12.6	1.0708	23	154.25	67.75
2	6.9	1.0853	22	173.25	72.25
3	24.6	1.0414	22	154.00	66.25
4	10.9	1.0751	26	184.75	72.25
5	27.8	1.0340	24	184.25	71.25
6	20.6	1.0502	24	210.25	74.75
7	19.0	1.0549	26	181.00	69.75
8	12.8	1.0704	25	176.00	72.50
9	5.1	1.0900	25	191.00	74.00
10	12.0	1.0722	23	198.25	73.50

- ✱ x : HIGHT, y: WEIGHT

Chicago social economic census

- ✱ The census included 77 communities in Chicago
- ✱ The census evaluated the average hardship index of the residents
- ✱ The census evaluated the following parameters for each community:
 - ✱ PERCENT_OF_HOUSING_CROWDED
 - ✱ PERCENT_HOUSEHOLD_BELOW_POVERTY
 - ✱ PERCENT_AGED_16p_UNEMPLOYED
 - ✱ PERCENT_AGED_25p_WITHOUT_HIGH_SCHOOL_DIPLOMA
 - ✱ PERCENT_AGED_UNDER_18_OR_OVER_64
 - ✱ PER_CAPITA_INCOME

*Given a new community and its parameters,
can you predict its average hardship index with all these parameters?*

The regression problem



Some terminology

- ✱ Suppose the dataset $\{(\mathbf{x}, y)\}$ consists of N labeled items (\mathbf{x}_i, y_i)
- ✱ If we represent the dataset as a table
 - ✱ The d columns representing $\{\mathbf{x}\}$ are called **explanatory variables** $\mathbf{x}^{(j)}$
 - ✱ The numerical column y is called the **dependent variable**

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1	3	0
2	3	2
3	6	5

Variables of the Chicago census

- [1] "PERCENT_OF_HOUSING_CROWDED"
- [2] "PERCENT_HOUSEHOLDS_BELOW_POVERTY"
- [3] "PERCENT_AGED_16p_UNEMPLOYED"
- [4] "PERCENT_AGED_25p_WITHOUT_HIGH_SCHOOL_DIPLOMA"
- [5] "PERCENT_AGED_UNDER_18_OR_OVER_64"
- [6] "PER_CAPITA_INCOME"
- [7] "HardshipIndex"

Which is the dependent variable in the census example?

- A. "PERCENT_OF_HOUSING_CROWDED"
- B. "PERCENT_AGED_25p_WITHOUT_HIGH_SCHOOL_DIPLOMA"
- C. "HardshipIndex"
- D. "PERCENT_AGED_UNDER_18_OR_OVER_64"

Linear model

- ✱ We begin by modeling y as a linear function of $\mathbf{x}^{(j)}$ plus randomness

$$y = \mathbf{x}^{(1)}\beta_1 + \mathbf{x}^{(2)}\beta_2 + \dots + \mathbf{x}^{(d)}\beta_d + \xi$$

Where ξ is a zero-mean random variable that represents model error

- ✱ In vector notation:

$$y = \mathbf{x}^T \boldsymbol{\beta} + \xi$$

Where $\boldsymbol{\beta}$ is the d -dimensional vector of coefficients that we train

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1	3	0
2	3	2
3	6	5

Each data item gives an equation

✿ The model: $y = \mathbf{x}^T \boldsymbol{\beta} + \xi = \mathbf{x}^{(1)} \beta_1 + \mathbf{x}^{(2)} \beta_2 + \xi$

Training data

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1	3	0
2	3	2
3	6	5

Which together form a matrix equation

✿ The model $y = \mathbf{x}^T \boldsymbol{\beta} + \xi = \mathbf{x}^{(1)} \beta_1 + \mathbf{x}^{(2)} \beta_2 + \xi$

Training data

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1	3	0
2	3	2
3	6	5

$$\begin{bmatrix} 0 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 1 & 3 \\ 2 & 3 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \xi_1 \\ \xi_2 \\ \xi_3 \end{bmatrix}$$

Which together form a matrix equation

✿ The model $y = \mathbf{x}^T \boldsymbol{\beta} + \xi = \mathbf{x}^{(1)} \beta_1 + \mathbf{x}^{(2)} \beta_2 + \xi$

Training data

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1	3	0
2	3	2
3	6	5

$$\begin{bmatrix} 0 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 1 & 3 \\ 2 & 3 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \xi_1 \\ \xi_2 \\ \xi_3 \end{bmatrix}$$

$$\mathbf{y} = \mathbf{X} \cdot \boldsymbol{\beta} + \mathbf{e}$$

Q. What's the dimension of matrix X?

A. $N \times d$

B. $d \times N$

C. $N \times N$

D. $d \times d$

Training the model is to choose β

- ✱ Given a training dataset $\{(\mathbf{x}, y)\}$, we want to fit a model $y = \mathbf{x}^T \beta + \xi$

- ✱ Define $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$ and $X = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}$ and $\mathbf{e} = \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_N \end{bmatrix}$

- ✱ To train the model, we need to choose β that makes \mathbf{e} small in the matrix equation $\mathbf{y} = X \cdot \beta + \mathbf{e}$

Training using least squares

- ✱ In the least squares method, we aim to **minimize** $\|\mathbf{e}\|^2$

$$\|\mathbf{e}\|^2 = \|\mathbf{y} - X\boldsymbol{\beta}\|^2 = (\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta})$$

- ✱ Differentiating with respect to $\boldsymbol{\beta}$ and setting to zero

$$X^T X \boldsymbol{\beta} - X^T \mathbf{y} = 0$$

- ✱ If $X^T X$ is invertible, the least squares estimate of the coefficient is:

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$$

Derivation of least square solution



Least square solution in the project



Convex set and convex function

- ✱ If a set is convex, any line connecting two points in the set is completely included in the set



(a)



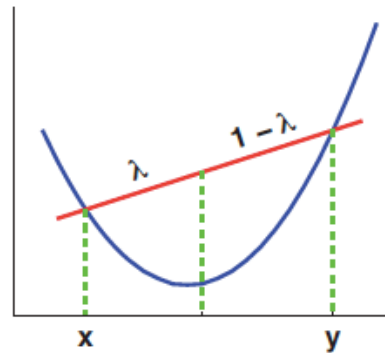
(b)

Figure 7.4 (a) Illustration of a convex set. (b) Illustration of a nonconvex set.

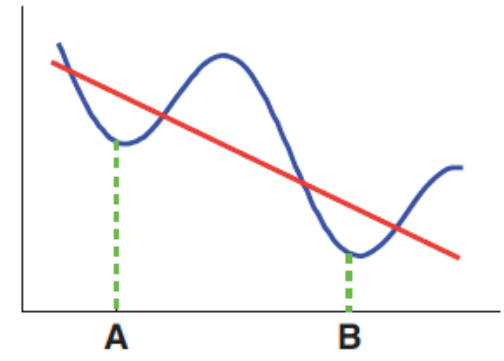
- ✱ A convex function: the area above the curve is convex

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$$

- ✱ The least square function is **convex**



(a)



(b)

What's the dimension of matrix $X^T X$?

- A. $N \times d$
- B. $d \times N$
- C. $N \times N$
- D. $d \times d$

Is this statement true?

If the matrix $\mathbf{X}^T\mathbf{X}$ does NOT have zero valued eigenvalues, it is invertible.

- A. TRUE
- B. FALSE

Training using least squares example

✿ Model: $y = \mathbf{x}^T \boldsymbol{\beta} + \xi = \mathbf{x}^{(1)} \beta_1 + \mathbf{x}^{(2)} \beta_2 + \xi$

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} = \begin{bmatrix} 2 \\ -\frac{1}{3} \end{bmatrix}$$

Training data

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1	3	0
2	3	2
3	6	5

$$\hat{\beta}_1 = 2$$
$$\hat{\beta}_2 = -\frac{1}{3}$$

Prediction

- ✱ If we train the model coefficients $\hat{\beta}$, we can predict y_0^p from \mathbf{x}_0

$$y_0^p = \mathbf{x}_0^T \hat{\beta}$$

- ✱ In the model $y = \mathbf{x}^{(1)}\beta_1 + \mathbf{x}^{(2)}\beta_2 + \xi$ with $\hat{\beta} = \begin{bmatrix} 2 \\ -\frac{1}{3} \end{bmatrix}$

- ✱ The prediction for $\mathbf{x}_0 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ is y_0^p

- ✱ The prediction for $\mathbf{x}_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ is y_0^p

A linear model with constant offset

✱ The problem with the model $y = \mathbf{x}^{(1)}\beta_1 + \mathbf{x}^{(2)}\beta_2 + \xi$

is:

✱ Let's add a constant offset β_0 to the model

$$y = \beta_0 + \mathbf{x}^{(1)}\beta_1 + \mathbf{x}^{(2)}\beta_2 + \xi$$

Training and prediction with constant offset

✱ The model $y = \beta_0 + \mathbf{x}^{(1)}\beta_1 + \mathbf{x}^{(2)}\beta_2 + \xi = \mathbf{x}^T\boldsymbol{\beta} + \xi$

✱ Training data:

$$\begin{bmatrix} 1 & x^{(1)} & x^{(2)} \end{bmatrix}$$

$\mathbf{1}$	$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1	1	3	0
1	2	3	2
1	3	6	5

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} = \begin{bmatrix} -3 \\ 2 \\ \frac{1}{3} \end{bmatrix}$$

✱ For $\mathbf{x}_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$y_0^p = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} -3 \\ 2 \\ \frac{1}{3} \end{bmatrix} = -3$$

Variance of the linear regression model

- ✱ The least squares estimate satisfies this property

$$\text{var}(\{y_i\}) = \text{var}(\{\mathbf{x}_i^T \hat{\boldsymbol{\beta}}\}) + \text{var}(\{\xi_i\})$$

- ✱ The random error is uncorrelated to the least square solution of linear combination of explanatory variables.

Variance of the linear regression model: proof

- ✱ The least squares estimate satisfies this property

$$\text{var}(\{y_i\}) = \text{var}(\{\mathbf{x}_i^T \hat{\boldsymbol{\beta}}\}) + \text{var}(\{\xi_i\})$$

Proof:

Evaluating models using R-squared

- ✱ The least squares estimate satisfies this property

$$\text{var}(\{y_i\}) = \text{var}(\{\mathbf{x}_i^T \hat{\boldsymbol{\beta}}\}) + \text{var}(\{\xi_i\})$$

- ✱ This property gives us an evaluation metric called R-squared

$$R^2 = \frac{\text{var}(\{\mathbf{x}_i^T \hat{\boldsymbol{\beta}}\})}{\text{var}(\{y_i\})}$$

- ✱ We have $0 \leq R^2 \leq 1$ with a larger value meaning a better fit.

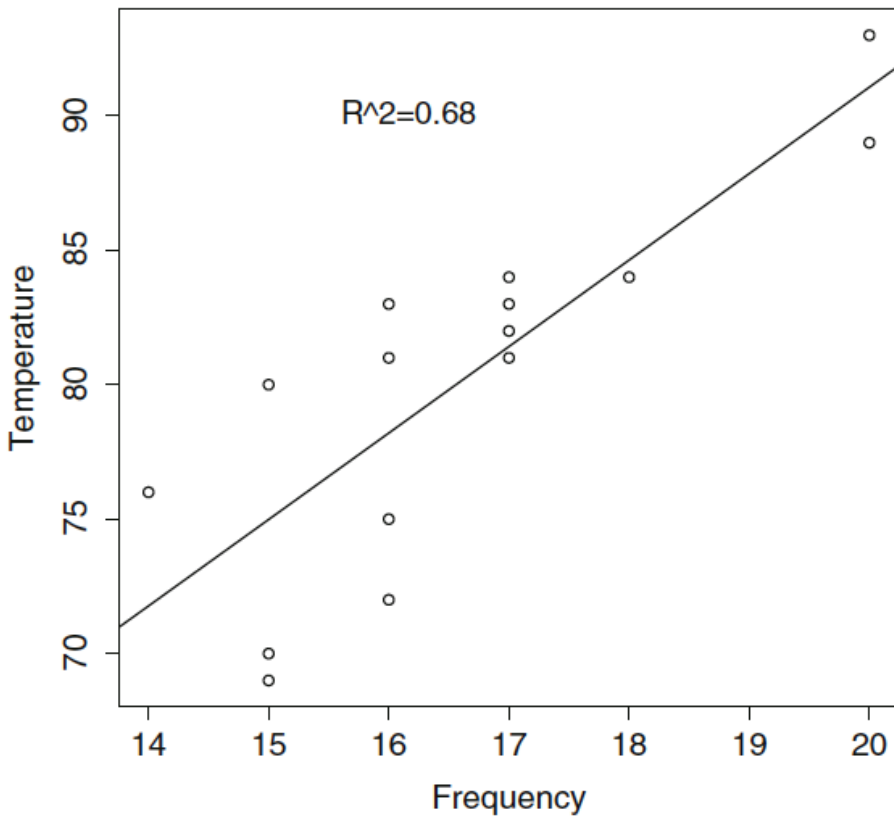
Q: What is R-squared if there is only one explanatory variable in the model?

Q: What is R-squared if there is only one explanatory variable in the model?

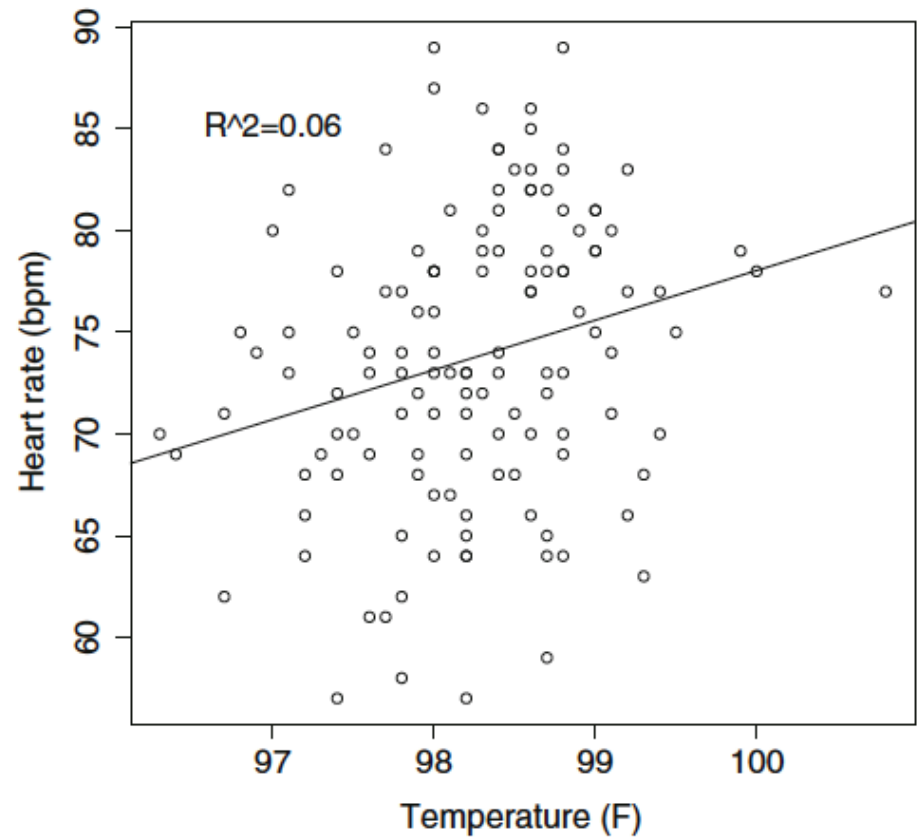
R-squared would be **the correlation coefficient squared** (textbook pgs 43-44)

R-squared examples

Chirp frequency vs temperature in crickets



Heart rate vs temperature in humans



Comparing our example models

$$y = \mathbf{x}^{(1)}\beta_1 + \mathbf{x}^{(2)}\beta_2 + \xi$$

$$y = \beta_0 + \mathbf{x}^{(1)}\beta_1 + \mathbf{x}^{(2)}\beta_2 + \xi$$

$$\hat{\beta} = \begin{bmatrix} 2 \\ -\frac{1}{3} \end{bmatrix}$$

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y	$\mathbf{x}^T \hat{\beta}$
1	3	0	1
2	3	2	3
3	6	5	4

$$\hat{\beta} = \begin{bmatrix} -3 \\ 2 \\ \frac{1}{3} \end{bmatrix}$$

1	$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y	$\mathbf{x}^T \hat{\beta}$
1	1	3	0	0
1	2	3	2	2
1	3	6	5	5

Linear regression model for the Chicago census data

Call:

```
lm(formula = HardshipIndex ~ ., data = dat)
```

Residuals:

Min	1Q	Median	3Q	Max
-15.7157	-1.9230	0.1301	1.9810	8.6719

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	105.1394	37.3622	2.814	0.006346	**
PERCENT_OF_HOUSING_CROWDED	0.7189	0.2753	2.612	0.011014	*
PERCENT_HOUSEHOLDS_BELOW_POVERTY	0.6665	0.0781	8.534	1.90e-12	***
PERCENT_AGED_16p_UNEMPLOYED	0.8023	0.1350	5.941	9.93e-08	***
PERCENT_AGED_25p_WITHOUT_HIGH_SCHOOL_DIPLOMA	0.7751	0.1063	7.293	3.64e-10	***
PERCENT_AGED_UNDER_18_OR_OVER_64	0.4807	0.1202	3.998	0.000156	***
PER_CAPITA_INCOME	-11.8819	3.1888	-3.726	0.000391	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

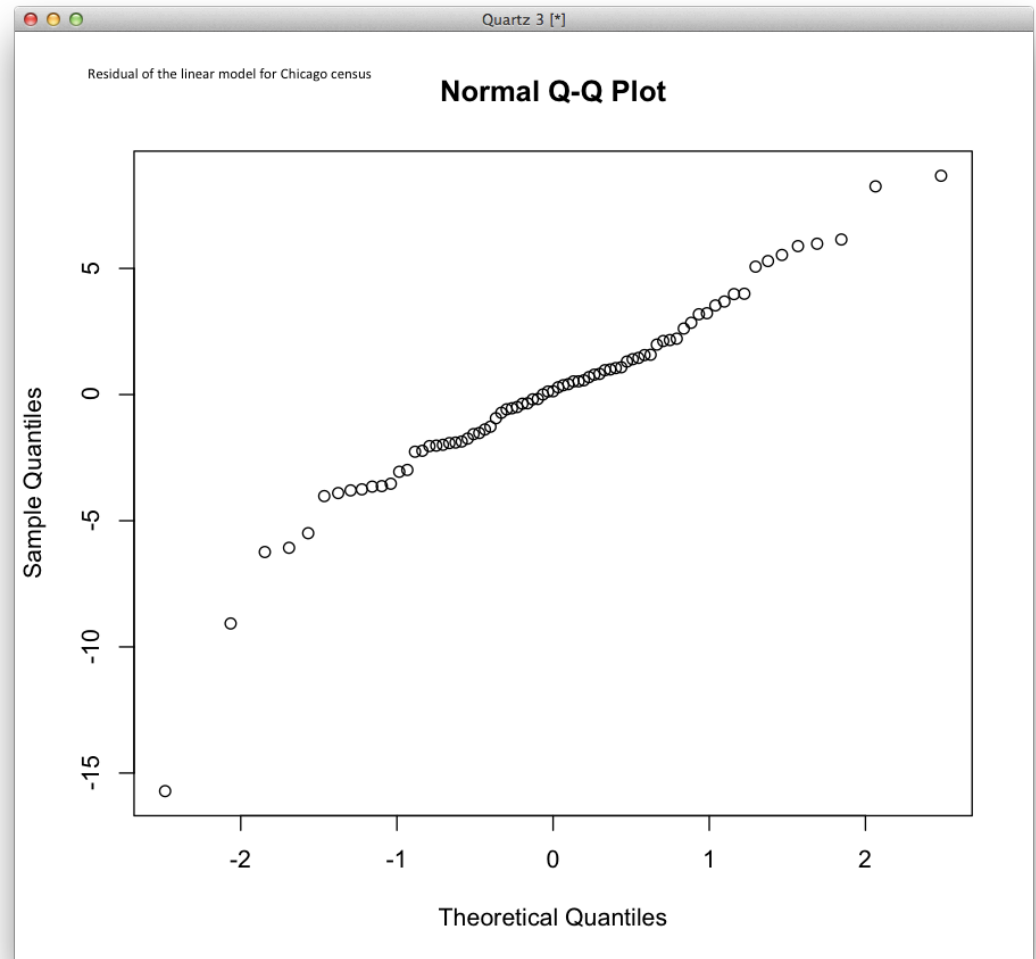
Residual standard error: 3.9 on 70 degrees of freedom

Multiple R-squared: 0.983, Adjusted R-squared: 0.9815

F-statistic: 673.9 on 6 and 70 DF, p-value: < 2.2e-16

Residual is normally distributed?

The Q-Q plot of the residuals is roughly normal



Prediction for another community

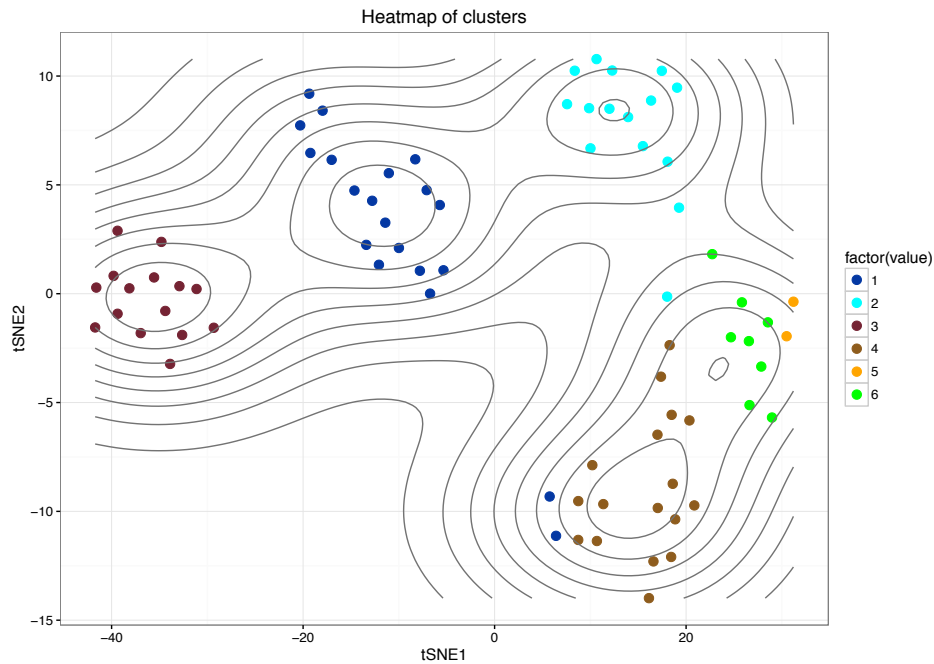
[1] "PERCENT_OF_HOUSING_CROWDED"	4.7
[2]"PERCENT_HOUSEHOLDS_BELOW_POVERTY"	19.7
[3] "PERCENT_AGED_16p_UNEMPLOYED"	12.9
[4]"PERCENT_AGED_25p_WITHOUT_HIGH_SCHOOL_DIPLOMA"	19.5
[5] "PERCENT_AGED_UNDER_18_OR_OVER_64"	33.5
[6]"PER_CAPITA_INCOME"	Log(28202)

Predicted hardship index: **41.46038**

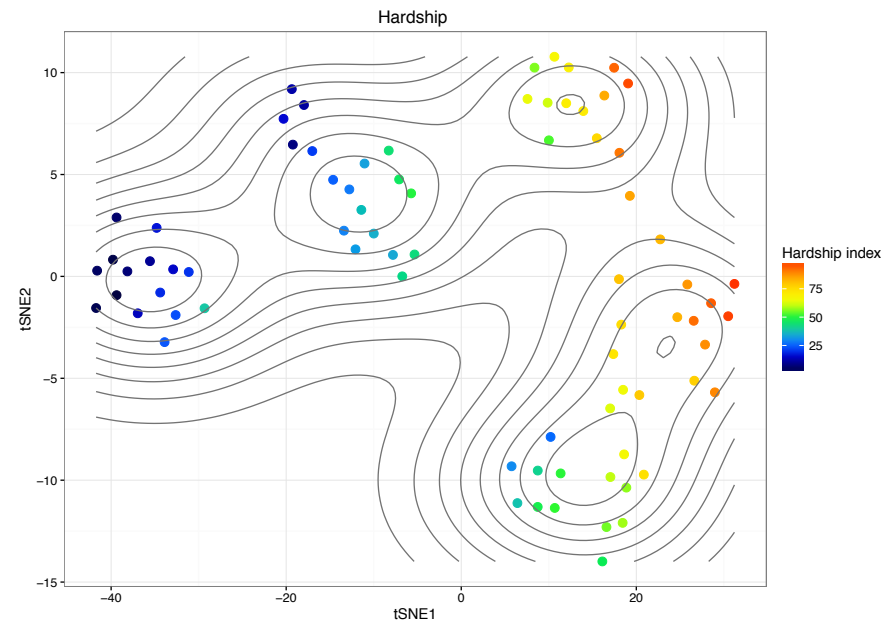
Note: maximum of hardship index in the training data is 98, minimum is 1

The clusters of the Chicago communities: clusters and hardship

Clusters of community



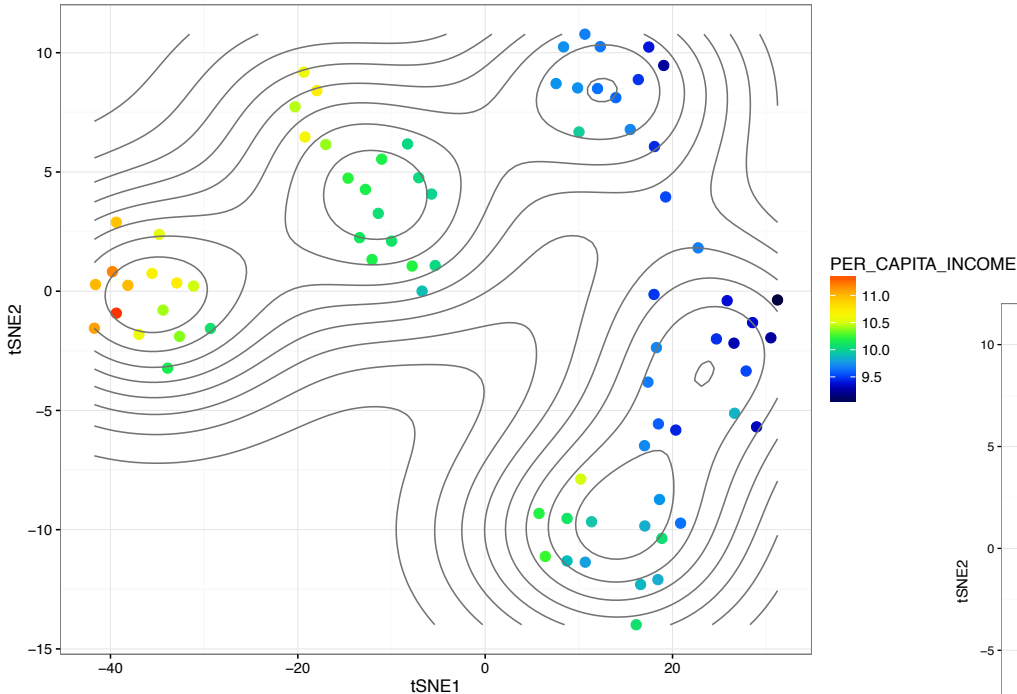
Hardship index of communities



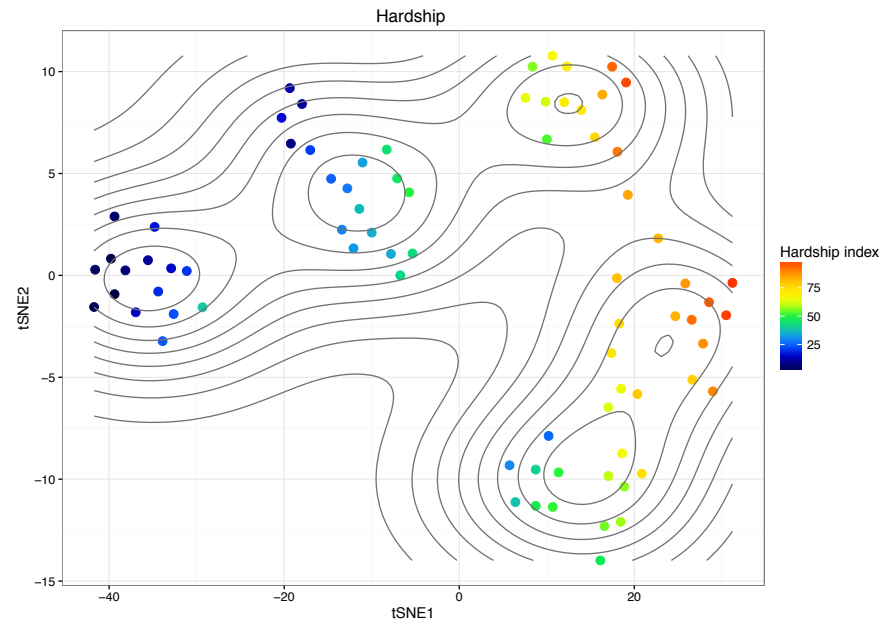
The clusters of the Chicago communities: per capital income and hardship

PER_CAPITAL_INCOME

Heatmap of PER_CAPITA_INCOME (log scale)

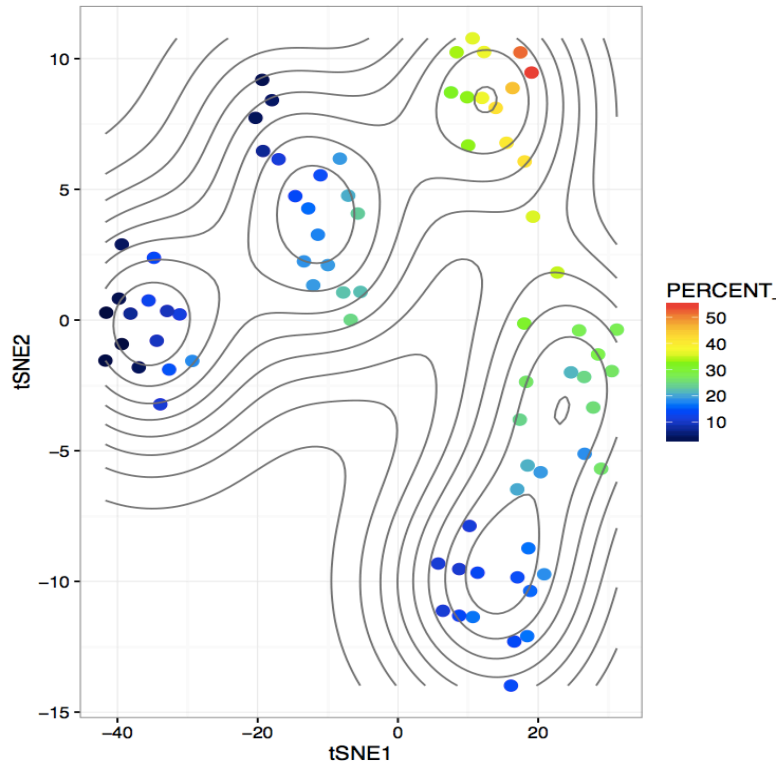


Hardship index of communities

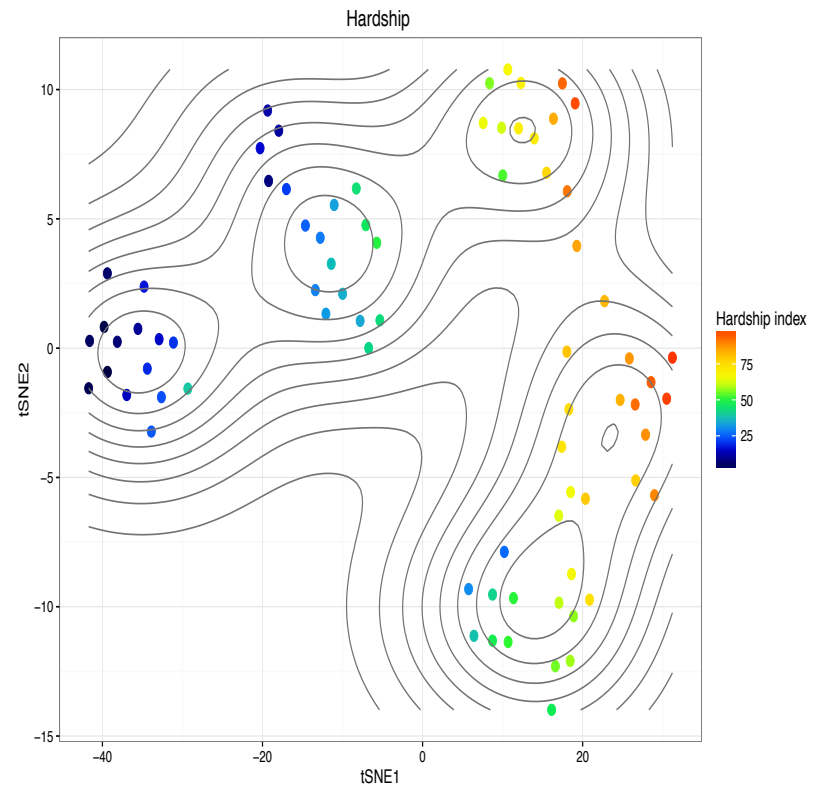


The clusters of the Chicago communities: without diploma and hardship

PERCENT_AGED_25p_WITHOUT_HIGH_SCHOOL_DIPLOMA



Hardship index of communities



Assignments

- ✱ Read Chapter 13 of the textbook
- ✱ Week 13 module
- ✱ Next time: More on linear regression

Additional References

- ✱ Robert V. Hogg, Elliot A. Tanis and Dale L. Zimmerman. “Probability and Statistical Inference”
- ✱ Kelvin Murphy, “Machine learning, A Probabilistic perspective”

See you next time

*See
You!*

