

IS590Labtime

topic: Expectation-maximization (EM) algorithm and  
Latent Dirichlet Allocation (LDA) algorithm

Date: 10/18

Time: 4:40-5:30pm

Instructor: Yingjun Guan

# A big picture of Machine Learning.

## Super VIP Cheatsheet: Machine Learning

Afshine AMIDI and Shervine AMIDI

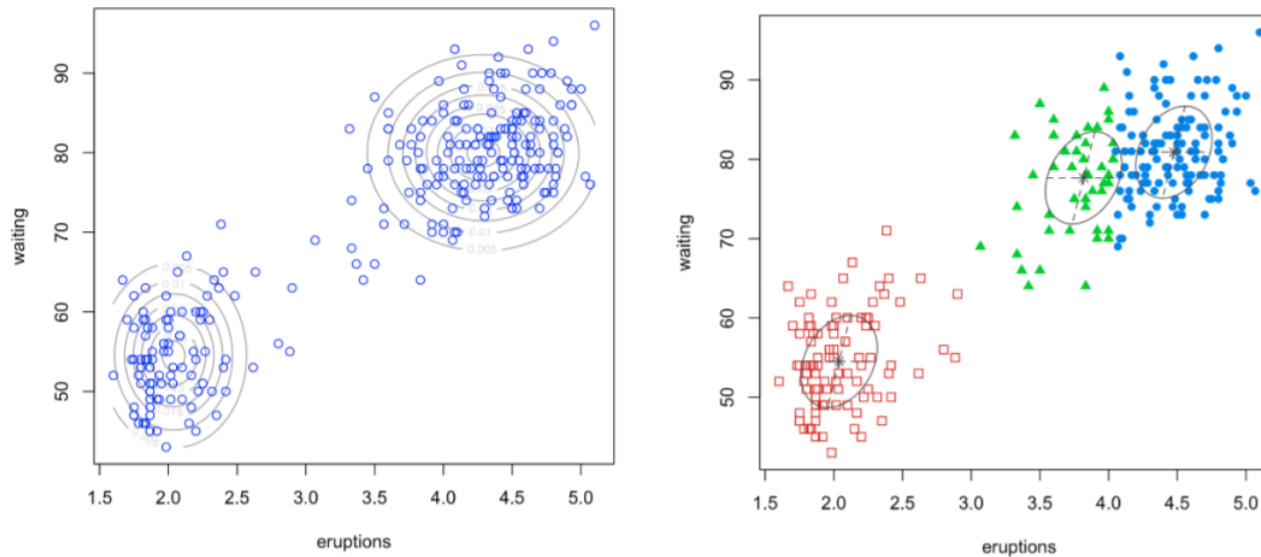
October 6, 2018

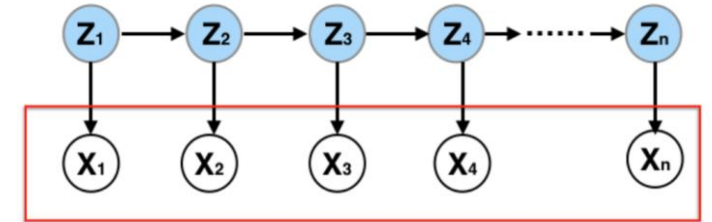
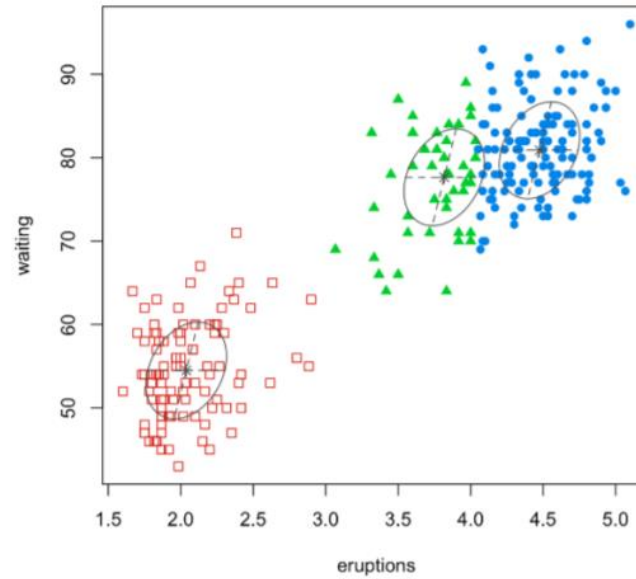
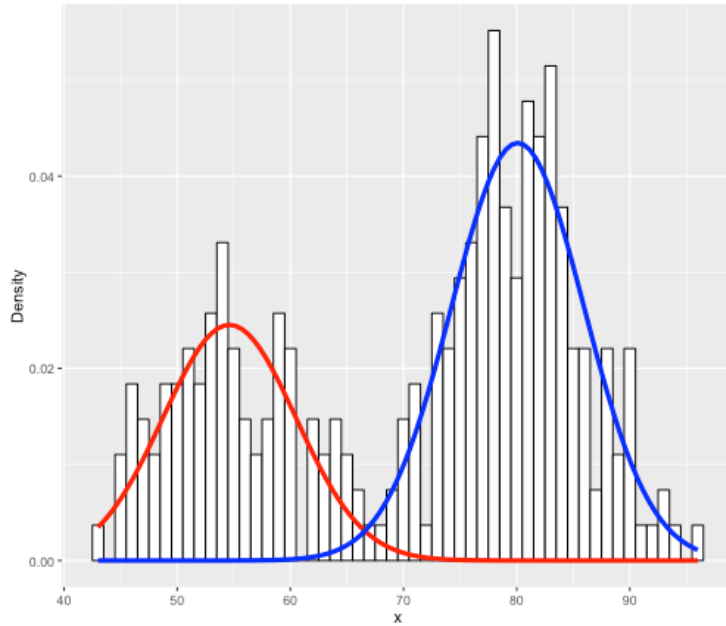
### Contents

<b>1 Supervised Learning</b>	<b>2</b>	<b>4 Machine Learning Tips and Tricks</b>	<b>10</b>
1.1 Introduction to Supervised Learning	2	4.1 Metrics	10
1.2 Notations and general concepts	2	4.1.1 Classification	10
1.3 Linear models	2	4.1.2 Regression	10
1.3.1 Linear regression	2	4.2 Model selection	11
1.3.2 Classification and logistic regression	3	4.3 Diagnostics	11
1.3.3 Generalized Linear Models	3		
1.4 Support Vector Machines	3	<b>5 Refreshers</b>	<b>12</b>
1.5 Generative Learning	4	5.1 Probabilities and Statistics	12
1.5.1 Gaussian Discriminant Analysis	4	5.1.1 Introduction to Probability and Combinatorics	12
1.5.2 Naive Bayes	4	5.1.2 Conditional Probability	12
1.6 Tree-based and ensemble methods	4	5.1.3 Random Variables	13
1.7 Other non-parametric approaches	4	5.1.4 Jointly Distributed Random Variables	13
1.8 Learning Theory	5	5.1.5 Parameter estimation	14
		5.2 Linear Algebra and Calculus	14
<b>2 Unsupervised Learning</b>	<b>6</b>	5.2.1 General notations	14
2.1 Introduction to Unsupervised Learning	6	5.2.2 Matrix operations	15
2.2 Clustering	6	5.2.3 Matrix properties	15
2.2.1 Expectation-Maximization	6	5.2.4 Matrix calculus	16
2.2.2 $k$ -means clustering	6		
2.2.3 Hierarchical clustering	6		
2.2.4 Clustering assessment metrics	6		
2.3 Dimension reduction	7		
2.3.1 Principal component analysis	7		
2.3.2 Independent component analysis	7		
<b>3 Deep Learning</b>	<b>8</b>		
3.1 Neural Networks	8		
3.2 Convolutional Neural Networks	8		
3.3 Recurrent Neural Networks	8		
3.4 Reinforcement Learning and Control	9		

# Today, we're talking about model-based clustering.

Model-based clustering refers to clustering a set of data points  $(x_1, \dots, x_n)$  by fitting a **mixture model** on this data set, where each cluster corresponds to a component of the mixture model.





Latent variables are hidden/unobserved data help to estimate the clustering, thus could be described in different ways.

---

# What is EM?

- The Expectation-Maximization (EM) algorithm is an iterative method that finds the MLE by enlarging the sample with **unobserved latent data**.

**Parameter  $\rightarrow$  Probability** [Expectation]

- E-step: Evaluate the posterior probability  $Q_i(z^{(i)})$  that each data point  $x^{(i)}$  came from a particular cluster  $z^{(i)}$  as follows:

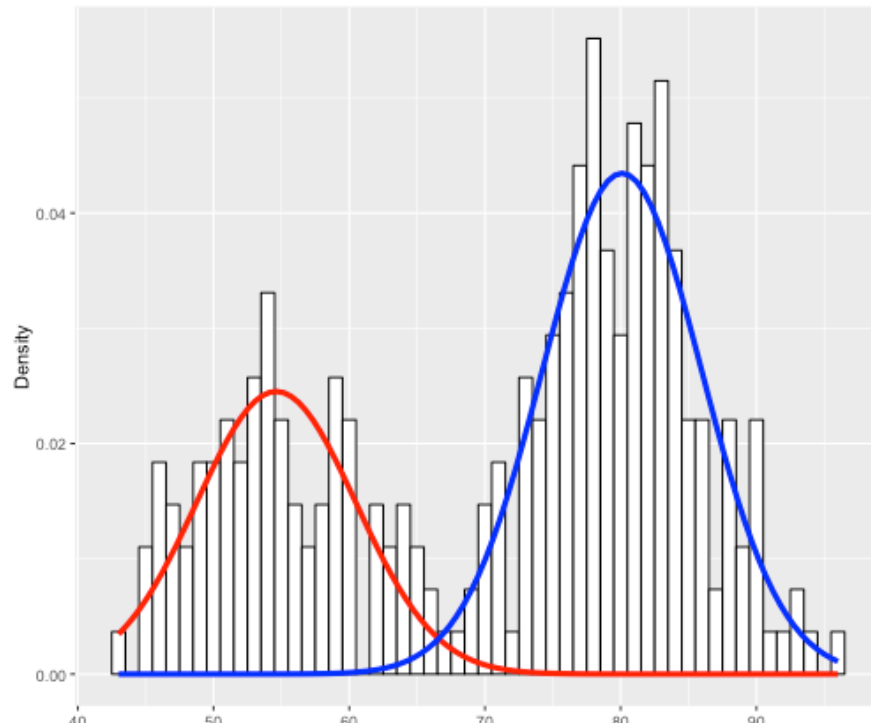
$$Q_i(z^{(i)}) = P(z^{(i)}|x^{(i)}; \theta)$$

- M-step: Use the posterior probabilities  $Q_i(z^{(i)})$  as cluster specific weights on data points  $x^{(i)}$  to separately re-estimate each cluster model as follows:

$$\theta_i = \operatorname{argmax}_{\theta} \sum_i \int_{z^{(i)}} Q_i(z^{(i)}) \log \left( \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) dz^{(i)}$$

**Probability  $\rightarrow$  Parameter** [Maximization]

# How to explain EM in equations (n-component gaussian mixture)



$$\log p(\mathbf{x}|\theta) = \sum_{i=1}^n \log \left[ \pi \phi_{\mu_1, \sigma_1^2}(x_i) + (1 - \pi) \phi_{\mu_2, \sigma_2^2}(x_i) \right].$$

- **E-step:** Let  $\theta_0$  denote the current value of  $\theta$ . Find  $p(\mathbf{Z}|\mathbf{x}, \theta_0)$ , the distribution of the latent variable  $\mathbf{Z}$  given the data  $\mathbf{x}$  and  $\theta_0$ , and then calculate the following expectation

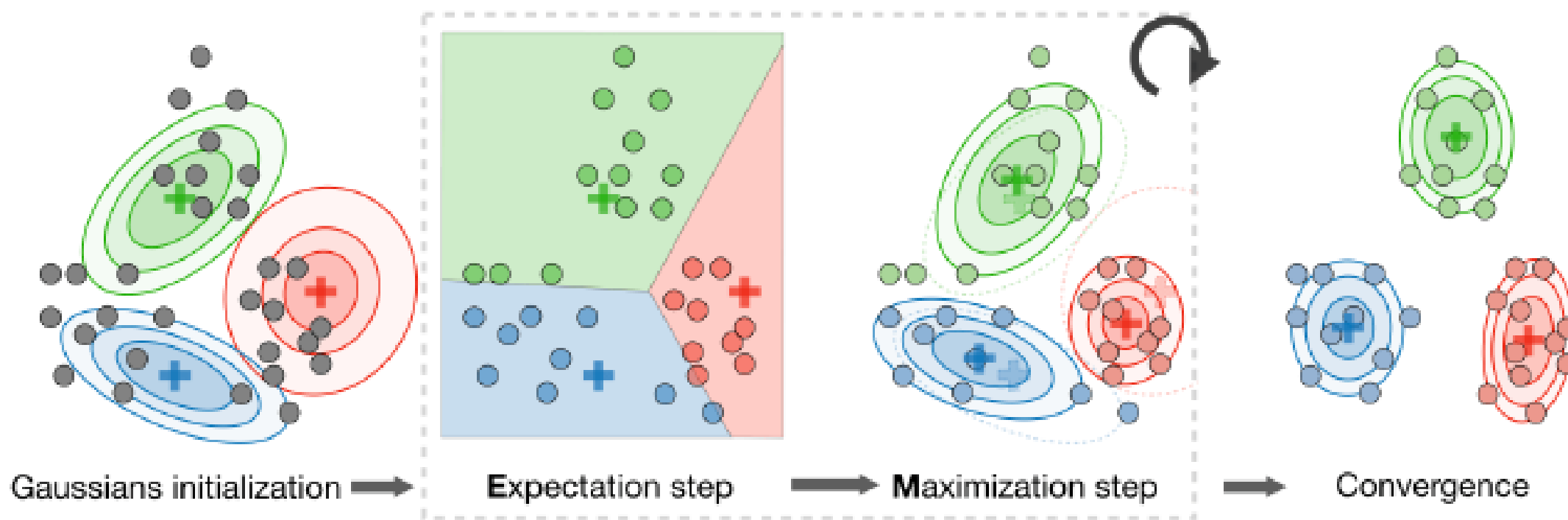
$$g(\theta) = \mathbb{E}_{\mathbf{Z}|\mathbf{x}, \theta_0} \log p(\mathbf{x}, \mathbf{Z}|\theta)$$

which is

$$\sum_{\mathbf{z}} p(\mathbf{Z} = \mathbf{z}|\mathbf{x}, \theta_0) \log p(\mathbf{x}, \mathbf{z}|\theta), \quad \text{or} \quad \int p(\mathbf{z}|\mathbf{x}, \theta_0) \log p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z}.$$

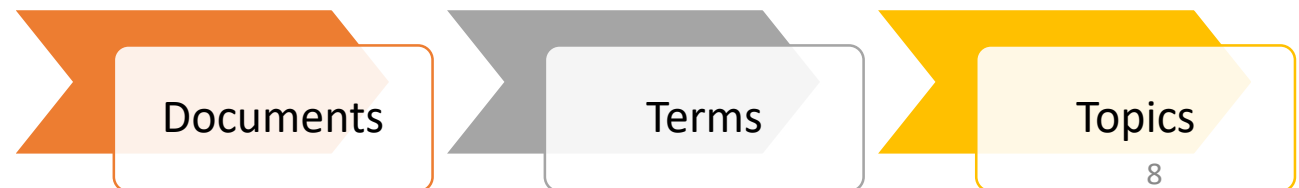
- **M-step:** Find  $\theta_1$  that maximizes  $g(\theta)$ .
- Replace  $\theta_0$  by  $\theta_1$  and repeat the above E and M steps until convergence.

# How to explain EM in figures. (n-component gaussian mixture)



# Hypotheses and targets of LDA (Latent Dirichlet Allocation)

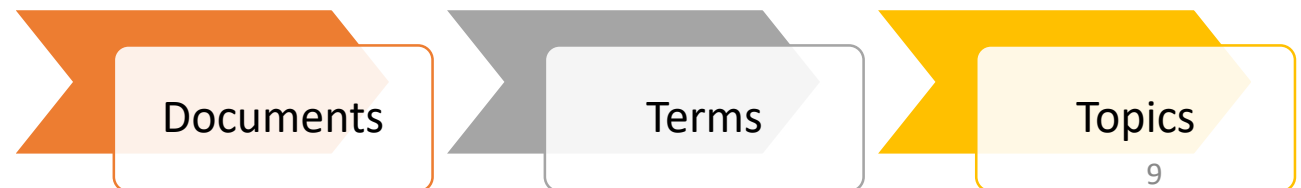
- We want to find themes (or topics) in documents
  - useful for e.g. search or browsing
- We don't want to do supervised topic classification
  - rather not fix topics in advance nor do manual annotation
- Need an approach which automatically teases out the topics
- This is essentially a clustering problem
  - can think of both words and documents as being clustered





# Key Assumptions behind the LDA topic model

- Documents exhibit multiple topics (but typically not many)
- LDA is a probabilistic model with a corresponding generative process
  - each document is assumed to be generated by this (simple) process
- A topic is a distribution over a fixed vocabulary
  - these topics are assumed to be generated first, before the documents
- Only the number of topics is specified in advance



Now, let's  
look at the  
generative  
process of  
LDA.

To generate a document:

1. Randomly choose a distribution over topics
  2. For each word in the document
    - a. randomly choose a topic from the distribution over topics
    - b. randomly choose a word from the corresponding topic (distribution over the vocabulary)
- Note that we need a distribution over a distribution (for step 1)
  - Note that words are generated independently of other words (unigram bag-of-words model)

Documents

Terms

Topics

## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

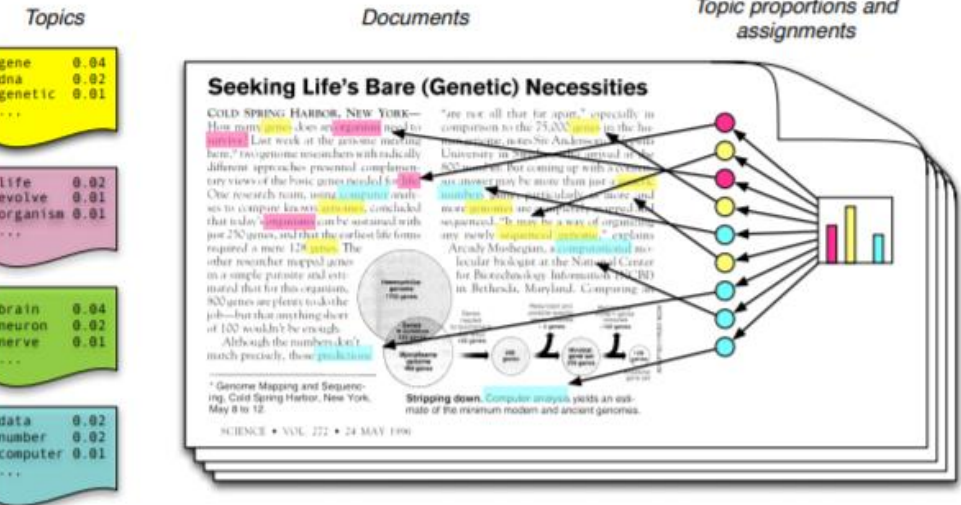
"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers game**, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains

Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

**Simple intuition:** Documents exhibit multiple topics.



- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those topics

Documents → terms → topics

# R-implementation

- Are you with me so far?
- Any questions are welcome.
- Thank you for the feedback last week.
  - Time
  - Hands-on practice
  - Topics
  - Pace/background

TF-IDF: one way of textual feature selection

# TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

Term frequency

Number of times term  $t$  appears in a doc,  $d$

Inverse document frequency

$$\log \frac{1 + n}{1 + \text{df}(d, t)}$$

# of documents  $n$

Document frequency of the term  $t$   $\text{df}(d, t)$

# Reference:

- Online tutorial.  
[https://www.cl.cam.ac.uk/teaching/1213/L101/clark\\_lectures/lect7.pdf](https://www.cl.cam.ac.uk/teaching/1213/L101/clark_lectures/lect7.pdf)
- David Blei's webpage is a good place to start
- Intro to EM on Youtube  
[https://www.youtube.com/watch?v=REypj2sy\\_5U](https://www.youtube.com/watch?v=REypj2sy_5U)
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan), 993-1022.
- <https://www.tidyttextmining.com/topicmodeling.html>

# Reference: (cont'd)

- Repository for today's lab:
  - Lecture slides
  - R file.
  - Other references.
- Online Access to today's lab.
  - 590DT student: log in through moodle directly.
  - Guest will be available through the following link:  
<https://us.bbcollab.com/guest/3DA671955178CE4AF10F31B05C983E27>
- Feedback for today's lab.
  - <https://goo.gl/forms/qpDmxEsYzUwbK7zx2>