

COMS 4721: Machine Learning for Data Science

Lecture 1, 1/17/2017

Prof. John Paisley

Department of Electrical Engineering
& Data Science Institute
Columbia University

OVERVIEW

This class will cover model-based techniques for extracting information from data with an end-task in mind. Such tasks include:

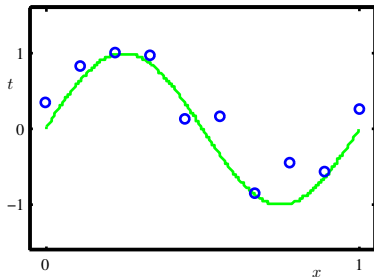
- ① ► predicting an unknown “output” given its corresponding “input”
- ② ► uncovering information within the data to better understand it
- ③ ► data-driven recommendation, grouping, classification, ranking, etc.

There are a few ways we can divide up the material as we go along, e.g.,

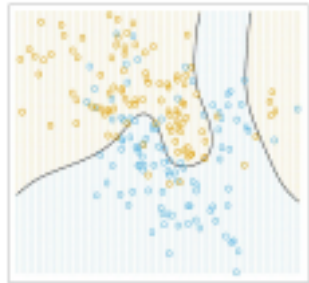
supervised learning		unsupervised learning
probabilistic models		non-probabilistic models
modeling approach		optimization techniques

We'll adopt the first method and work in the second two along the way.

OVERVIEW: SUPERVISED LEARNING



(a) Regression



(b) Classification

Regression: Using set of inputs, predict real-valued output.

Classification: Using set of inputs, predict a discrete label (aka class).

EXAMPLE CLASSIFICATION PROBLEM

Given a set of inputs characterizing an item, assign it a label.

Is this spam?

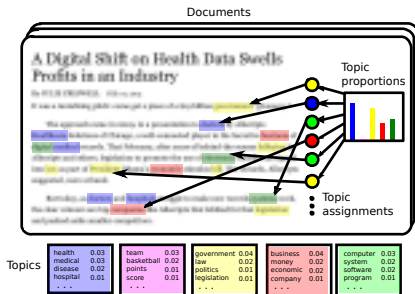
hi everyone,

i saw that close to my hotel there is a pub with bowling
(it's on market between 9th and 10th avenue). meet
there at 8:30?

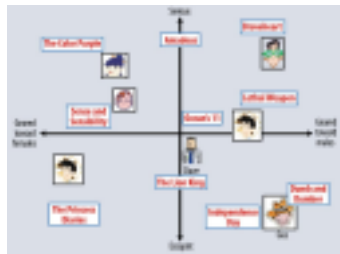
What about this?

Enter for a chance to win a trip to Universal Orlando to
celebrate the arrival of Dr. Seuss's The Lorax on Movies
On Demand on August 21st! [Click here now!](#)

OVERVIEW: UNSUPERVISED LEARNING



(c) topic modeling



(d) recommendations¹

With unsupervised learning our goal is often to uncover structure in the data. This helps with predictions, recommendations, efficient data exploration.

¹ Figure from Koren, Y., Robert B., and Volinsky, C.. "Matrix factorization techniques for recommender systems." Computer 42.8 (2009): 30-37.

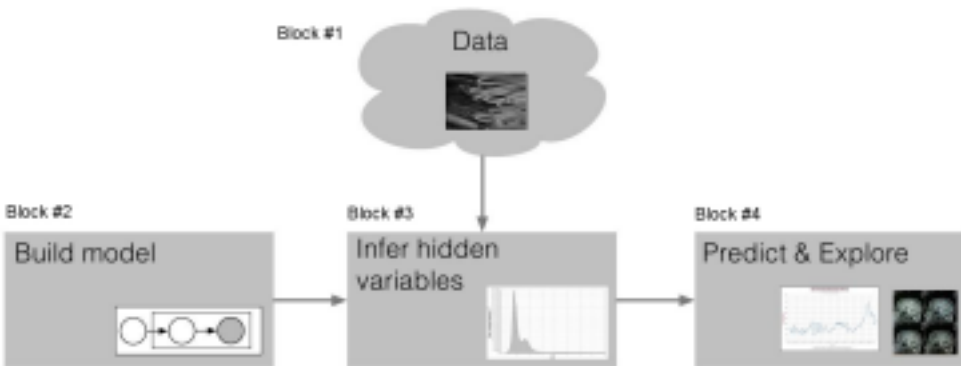
EXAMPLE UNSUPERVISED PROBLEM

Goal: Learn the dominant topics from a set of news articles.

The New York Times

music band songs rock album jazz pop song single night	book life novel story books man stories love children family	art museum show exhibition artist artists paintings painting century works	game knicks nets points team season play games night coach	show film television movie series says life man character know
theater play production show stage street broadway director musical directed	clinton bush campaign core political republican cole presidential senator house	stock market percent fund investors funds companies stocks investment trading	restaurant sauce menu food dishes street dining dinner chicken served	budget talk governor county mayor billion taxes plan legislature fiscal

DATA MODELING

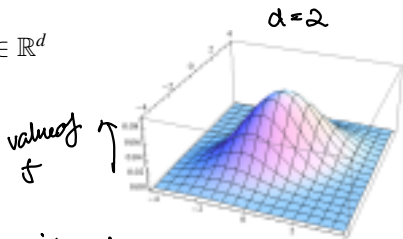


- ▶ Supervised vs. unsupervised: Blocks #1 and #4
- ▶ Probabilistic vs. non-probabilistic: Primarily Block #2 (Some Block #3)
- ▶ Model development (Block #2) vs. Optimization techniques (Block #3)

GAUSSIAN DISTRIBUTION (MULTIVARIATE)

Gaussian density in d dimensions

- ▶ Block #1: Data x_1, \dots, x_n . Each $x_i \in \mathbb{R}^d$
- ▶ Block #2: An i.i.d. Gaussian model
- ▶ Block #3: Maximum likelihood
- ▶ Block #4: Leave undefined



The density function is

Probability of x given parameters μ, Σ

$$p(x|\mu, \Sigma) := \frac{1}{(2\pi)^{\frac{d}{2}} \sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2} \underbrace{(x-\mu)^T \Sigma^{-1} (x-\mu)}_{\text{quadratic function}}\right)$$

positive definite

i/p vector

non-negative, integrator 1.

We evaluate this function and that gives a density at x using a gaussian with mean μ & covariance Σ

The central moments are:

$$\mathbb{E}[x] = \int_{\mathbb{R}^d} x p(x|\mu, \Sigma) dx = \mu,$$

subtract mean of x from x before.

and then evaluate expectation of this outer product

equal to sigma

$$\text{Cov}(x) = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^T] = \mathbb{E}[xx^T] - \mathbb{E}[x]\mathbb{E}[x]^T = \Sigma.$$

↳ spread of distribution in various dimensions, as well as the correlation between those dimensions

BLOCK #2: A PROBABILISTIC MODEL

Probabilistic Models

→ simply a set of probability distributions on our data.
We have data X , we have to define some probability distribution on X .
That distribution is going to take parameters

- ▶ A probabilistic model is a set of probability distributions, $p(x|\theta)$.
- ▶ We pick the distribution family $p(\cdot)$, but don't know the parameter θ .

No matter what value we set θ to, we are always working within the same distribution family.

Example: Model data with a Gaussian distribution $p(x|\theta)$, $\theta = \{\mu, \Sigma\}$.

(often made \downarrow)

this density is a multi-variate Gaussian & takes parameters μ & Σ .
For all parameters μ & Σ , we are working with the same family.

The i.i.d. assumption

Assume data is independent and identically distributed (iid). This is written every single data point is simply independent of every other and it has the same distribution family.

$$x_i \stackrel{iid}{\sim} p(x|\theta), \quad i = 1, \dots, n.$$

Writing the density as $p(x|\theta)$, then the joint density decomposes as
This assumption allows us to write our joint density this way \rightarrow

$$p(x_1, \dots, x_n|\theta) = \prod_{i=1}^n p(x_i|\theta).$$

Probability of n d -dimensional vectors coming from a distribution with parameter θ

Joint density (probability) of all n observations is simply the product of the prob. of each of them individually.

BLOCK #3: MAXIMUM LIKELIHOOD ESTIMATION

(define an objective function).

data & unknown model parameters

Maximum Likelihood approach (common approach with probabilistic models).

We now need to find θ . Maximum likelihood seeks the value of θ that maximizes the likelihood function:

For a particular value of θ , it tells us what is the likelihood of data coming from that model. $\hat{\theta}_{ML} := \arg \max_{\theta} p(x_1, \dots, x_n | \theta)$, which is the prob. of observations that we have given that particular value of θ .

[Find θ that maximizes joint likelihood function.]

uprob. distribution on our data [joint likelihood]

This value best explains the data according to the chosen distribution family. we want to find the value of θ that says this dataset is most probable from a m.v. distribution having this parameterization.

[In a sense we are fitting the data. $\hat{\theta}_{ML} \rightarrow$ maximum likelihood estimate]

Maximum Likelihood equation

The analytic criterion for this maximum likelihood estimator is:



$$\nabla_{\theta} \prod_{i=1}^n p(x_i | \theta) = 0.$$

Maximum of a function is the pt at which gradient = 0.

Simply put, the maximum is at a peak. There is no "upward" direction.

BLOCK #3: LOGARITHM TRICK *(This technique is problem specific)*

(When using prob. distributions in an iid setting.)

Logarithm trick

(log doesn't change location of maximum & minimum)

Calculating $\nabla_{\theta} \prod_{i=1}^n p(x_i|\theta)$ can be complicated. We use the fact that the logarithm is monotonically increasing on \mathbb{R}_+ , and the equality

$$\ln\left(\prod_i f_i\right) = \sum_i \ln(f_i).$$

product *sum*

Consequence: Taking the logarithm does not change the *location* of a maximum or minimum:

$$\max_y \ln g(y) \neq \max_y g(y)$$

($g \neq \ln g$)
The *value* changes. \nearrow

$$\arg \max_y \ln g(y) = \arg \max_y g(y)$$

The location does not change.

\downarrow
which we care about

BLOCK #3: ANALYTIC MAXIMUM LIKELIHOOD

Maximum likelihood and the logarithm trick

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} \overset{\text{iid assumption})}{\prod_{i=1}^n} p(x_i|\theta) = \arg \max_{\theta} \overset{\text{product}}{\ln\left(\prod_{i=1}^n p(x_i|\theta)\right)} = \arg \max_{\theta} \overset{\text{sum}}{\sum_{i=1}^n \ln p(x_i|\theta)}$$

To then solve for $\hat{\theta}_{\text{ML}}$, find

$$\nabla_{\theta} \sum_{i=1}^n \ln p(x_i|\theta) = \sum_{i=1}^n \nabla_{\theta} \ln p(x_i|\theta) = 0.$$

[we work with this function.]
[gradient at $\theta=0$]
Gradient of each individual x_i likelihood over each data point

Depending on the choice of the model, we will be able to solve this

1. analytically (via a simple set of equations) \rightarrow in this case
 2. numerically (via an iterative algorithm using different equations)
 3. approximately (typically when #2 converges to a local optimal solution)
- in complicated models we can't take the derivative & set it = 0.*

EXAMPLE: MULTIVARIATE GAUSSIAN MLE

Block #2: Multivariate Gaussian data model

Model: Set of all Gaussians on \mathbb{R}^d with unknown mean $\mu \in \mathbb{R}^d$ and covariance $\Sigma \in \mathbb{S}_{++}^d$ (positive definite $d \times d$ matrix).

We assume that x_1, \dots, x_n are i.i.d. $p(x|\mu, \Sigma)$, written $x_i \stackrel{iid}{\sim} p(x|\mu, \Sigma)$.

Block #3: Maximum likelihood solution

We have to solve the equation

$$\sum_{i=1}^n \nabla_{(\mu, \Sigma)} \ln p(x_i | \mu, \Sigma) = 0$$

for μ and Σ . (Try doing this without the log to appreciate it's usefulness.)

in case it is maximizing the f^ because we want it to be more probable.*

EXAMPLE: GAUSSIAN MEAN MLE

First take the gradient with respect to μ .

$$0 = \nabla_{\mu} \sum_{i=1}^n \ln \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right)$$

multi-variate Gaussian evaluated at each observation.

$$= \nabla_{\mu} \sum_{i=1}^n -\frac{1}{2} \ln(2\pi)^d |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

[Sum of log of each of those 2 terms]

$$0 = -\frac{1}{2} \sum_{i=1}^n \nabla_{\mu} \left(x_i^T \Sigma^{-1} x_i - 2\mu^T \Sigma^{-1} x_i + \mu^T \Sigma^{-1} \mu \right) = -\Sigma^{-1} \sum_{i=1}^n (x_i - \mu)$$

set = $\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$ vector (since it's not a scalar)

Since Σ is positive definite, the only solution is

$$\sum_{i=1}^n (x_i - \mu) = 0 \quad \Rightarrow \quad \hat{\mu}_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n x_i$$

Gradient of the log of joint likelihood with respect to mean vector μ .

Since this solution is independent of Σ , it doesn't depend on $\hat{\Sigma}_{\text{ML}}$.

EXAMPLE: GAUSSIAN COVARIANCE MLE

Now take the gradient with respect to Σ . ↗ (d x d matrix)

↖ matrix of 0s
d x d

$$\begin{aligned} 0 &= \nabla_{\Sigma} \sum_{i=1}^n -\frac{1}{2} \ln(2\pi)^d |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \\ &= -\frac{n}{2} \nabla_{\Sigma} \ln |\Sigma| - \frac{1}{2} \nabla_{\Sigma} \text{trace} \left(\Sigma^{-1} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T \right) \\ &= -\frac{n}{2} \Sigma^{-1} + \frac{1}{2} \Sigma^{-2} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T \end{aligned}$$

(function of μ & Σ)
↖ equating to matrix of 0s

Solving for Σ and plugging in $\mu = \hat{\mu}_{\text{ML}}$,

$$\hat{\Sigma}_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{\text{ML}})(x_i - \hat{\mu}_{\text{ML}})^T.$$

↳ empirical avg. of our outer products
of our data where we have subtracted
off max. likelihood mean.

EXAMPLE: GAUSSIAN MLE (SUMMARY)

So if we have data x_1, \dots, x_n in \mathbb{R}^d that we hypothesize is i.i.d. Gaussian, the maximum likelihood values of the mean and covariance matrix are

$$\hat{\mu}_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\Sigma}_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{\text{ML}})(x_i - \hat{\mu}_{\text{ML}})^T.$$

(empirical covariance of data where we subtract off max.-likelihood mean.)

Are we done? There are many assumptions/issues with this approach that makes finding the “best” parameter values not a complete victory.

- ▶ We made a model assumption (multivariate Gaussian).
- ▶ We made an i.i.d. assumption. (sequential information may be present)
- ▶ We assumed that maximizing the likelihood is the best thing to do. → why would it not be?

Comment: We often use θ_{ML} to make predictions about x_{new} (Block #4).

How does θ_{ML} generalize to x_{new} ?

If $x_{1:n}$ don't “capture the space” well, θ_{ML} can overfit the data.

The problem was solved for a function that we chose.