# COMS 4721: Machine Learning for Data Science
## Lecture 1, 1/17/2017

Prof. John Paisley

Department of Electrical Engineering
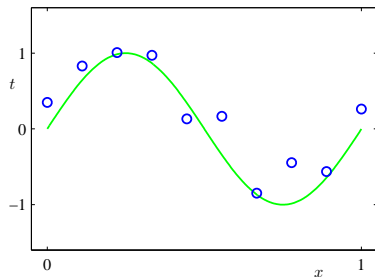& Data Science Institute

Columbia University

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

- predicting an unknown "output" given its corresponding "input"
- uncovering information within the data to better understand it
- data-driven recommendation, grouping, classification, ranking, etc.

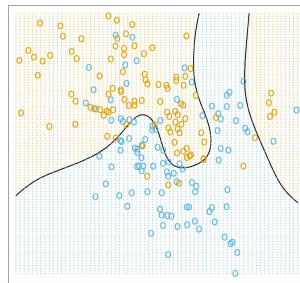There are a few ways we can divide up the material as we go along, e.g.,

| 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.

(a) Regression

(b) Classification

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

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

Given a set of inputs characterizing an item, assign it a label.
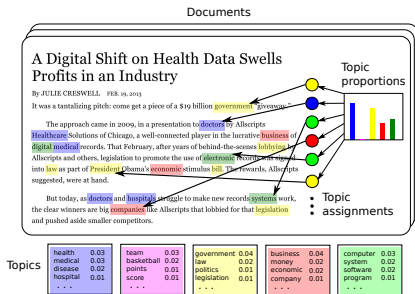
## Is this spam?

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
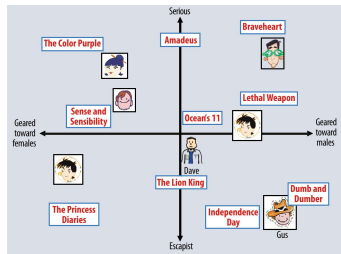there at 8:30?

## 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!

(c) topic modeling



(d) recommendations[1]

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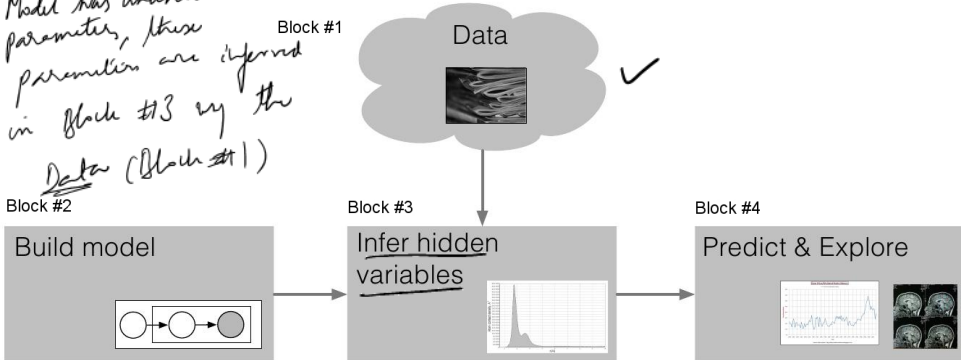
**Goal**: Learn the dominant topics from a set of news articles.

*The New York Times*

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

Model has unknown parameters, these parameters are inferred in Block #3 by the Data (Block #1)

- ▶ Supervised vs. unsupervised: Blocks #1 and #4

- ▶ Probabilistic vs. non-probabilistic: Primarily Block #2 (Some Block #3)
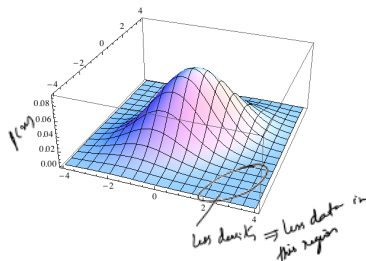
  Some algorithmic techniques → are purely motivated by fact that we are doing a probabilistic model.

  However a lot of techniques for non-prob models apply to prob model.

- ▶ Model development (Block #2) vs. Optimization techniques (Block #3)

## Gaussian density in $d$ dimensions

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

less density ⇒ less data in this region

The density function is

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

The central moments are:

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

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

### Probabilistic Models

- 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 $\theta$.

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

$\mu, \Sigma$ are inferred in Block #3

### The i.i.d. assumption

Assume data is *independent and identically distributed (iid)*. This is written

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

training data is produced independent y other samples

Writing the density as $p(x|\theta)$, then the *joint* density decomposes as

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

$\hookrightarrow$ Objective function: f (data & unknown parameter)

## Maximum Likelihood approach

We now need to find $\theta$. *Maximum likelihood* seeks the value of $\theta$ that maximizes the likelihood function:

$$\hat{\theta}_{\text{ML}} := \arg\max_\theta p(x_1, \ldots, x_n | \theta),$$

$\mapsto$ Ⓐ

This value best explains the data according to the chosen distribution family.

## Maximum Likelihood equation

The analytic criterion for this maximum likelihood estimator is:

$$\nabla_0 f(x|0) = 0$$

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

Max$^m$ of a function (i.e Ⓐ) is where gradient = 0.

$\mapsto$ Ⓐ

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

### Logarithm trick

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

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

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

$$\max_y \ln g(y) \neq \max_y g(y) \qquad \text{The } value \text{ changes.}$$

$$\boxed{\arg\max_y \ln g(y) = \arg\max_y g(y)} \qquad \text{The } location \text{ does not change.}$$

(true for positive function $g(y)$)

## Maximum likelihood and the logarithm trick

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} \prod_{i=1}^{n} p(x_i|\theta) = \arg \max_{\theta} \ln\Big(\prod_{i=1}^{n} p(x_i|\theta)\Big) = \arg \max_{\theta} \sum_{i=1}^{n} \ln p(x_i|\theta)$$

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

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

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

1. analytically (via a simple set of equations) $\smile$
2. numerically (via an iterative algorithm using different equations) $\smile$
3. approximately (typically when #2 converges to a local optimal solution)

$\rightarrow$ covariance matrix is always +ive
semi-definite.

## 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 $x_1, \ldots, x_n$ are i.i.d. $p(x|\mu, \Sigma)$, written $x_i \overset{iid}{\sim} p(x|\mu, \Sigma)$.

## 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 $\mu$ and $\Sigma$. (Try doing this without the log to appreciate it's usefulness.)

First take the gradient with respect to $\mu$.

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

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

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

Since $\Sigma$ is positive definite, the only solution is  *Generally $\Sigma$ covariance matrices are +ive semi-definite i.e $\geq 0$*

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

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

Now take the gradient with respect to $\Sigma$.

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

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

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

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

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

**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). $Block \#2$

▶ We made an i.i.d. assumption.

▶ We assumed that maximizing the likelihood is the best thing to do.

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

   How does $\theta_{\text{ML}}$ generalize to $x_{new}$? $\quad \rightarrow$ - MLE fails here.

   If $x_{1:n}$ don't "capture the space" well, $\theta_{\text{ML}}$ can *overfit* the data.

$\longrightarrow$ In block #3, we chose an objective function (i.e likelihood) to model the data. Later we maximized it cause we chose likelihood. Choosing likelihood they was an assumption that it is best thing to do.