

Generative models

Coin with some intrinsic probability P_H

and $P_T = 1 - P_H$

H T H H T H T H H T

Document classification

Categories - Spam / Not Spam
 P_S $1 - P_S$

Set of words $P(w_i | \text{spam})$ $P(w_i | \text{not spam})$

Generate a document

1. Choose Spam / not spam P_S
2. Depending on choice, for each w_i , include w_i with prob. $P(w_i | c)$

Naive model - "bag of words"

Bag = multiset = set with a count for each value

Generalize Spam/Not Spam to Topic

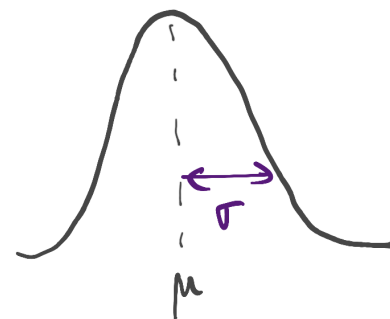
Modelling

Google News - automatically assign new stories to categories

What if some of the generating parameters are hidden?

Gaussian

$$-\frac{1}{2\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Given samples x_1, x_2, \dots, x_N

estimate σ, μ

Simpler: σ fixed, estimate μ

Given a hypothesis for $\mu - \hat{\mu}$

$$P(x_1) - \text{prob} \quad \frac{1}{\sigma^2} e^{-(x_1 - \hat{\mu})^2}$$

$$P(x_2) \dots$$

\vdots

$$P(x_n)$$

Choose $\hat{\mu}$ to maximize $\prod_i P(x_i | \hat{\mu})$

$$\text{Likelihood} = L(\theta, \mu)$$

Maximizing $L(\theta, \mu)$ is messy

Instead take logs

$$\ln(L(\theta, \mu)) = \ln \prod_{i=1}^k P(x_i | \mu)$$

$$= \sum_{i=1}^k \ln P(x_i | \mu)$$

Log

likelihood

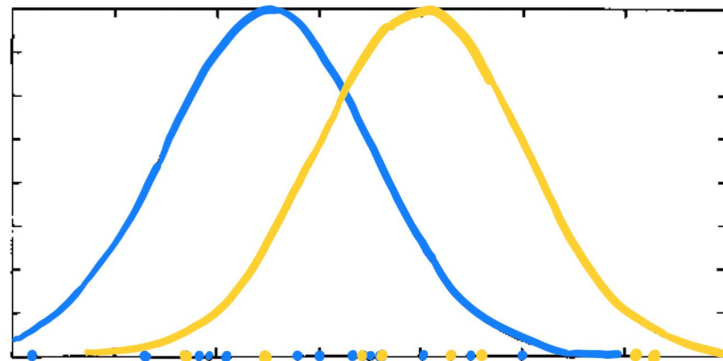
$$\ln \frac{1}{\sigma^2} e^{-(x_i - \mu)^2}$$

$$\text{Maximize } -(x_i - \mu)^2$$

Best choice for μ is sample

$$\text{mean} = \frac{1}{k} \sum_{i=1}^k x_i$$

Mixture of 2 Gaussians



Know colour of points - apply sample mean separately to blue/yellow points

If we don't know?

Generative model

Pick blue or yellow distribution (equally likely - uniform)

Generate x_i according to chosen distribution

H T T T H H T H T H
H H H H T H H H H H
H T H H H H H T H H
H T H T T T H H T T
T H H H T H H H T H

2 biased coins

Experiment

Choose one coin

(uniform) Repeat
K
times

Toss it M times

	Coin A	Coin B
H T T T H H T H T H		5 H, 5 T
H H H H T H H H H H	9 H, 1 T	
H T H H H H H T H H	8 H, 2 T	
H T H T T T H H T T		4 H, 6 T
T H H H T H H H T H	7 H, 3 T	

24 H 6 T

$P_H = 0.8$

9 H 11 T

$P_H = 0.45$

How to proceed if info about coins is missing

Initially Guess P_1, P_2
 P_H red P_H blue

For each row of 10 tosses - 3H, 7T

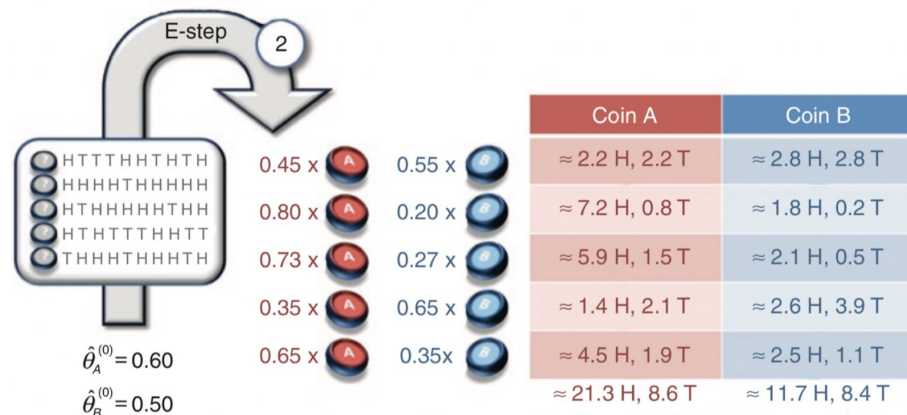
$P_1^3 (1-P_1)^7$ - red coin P_R

$P_2^3 (1-P_2)^7$ - blue coin P_B

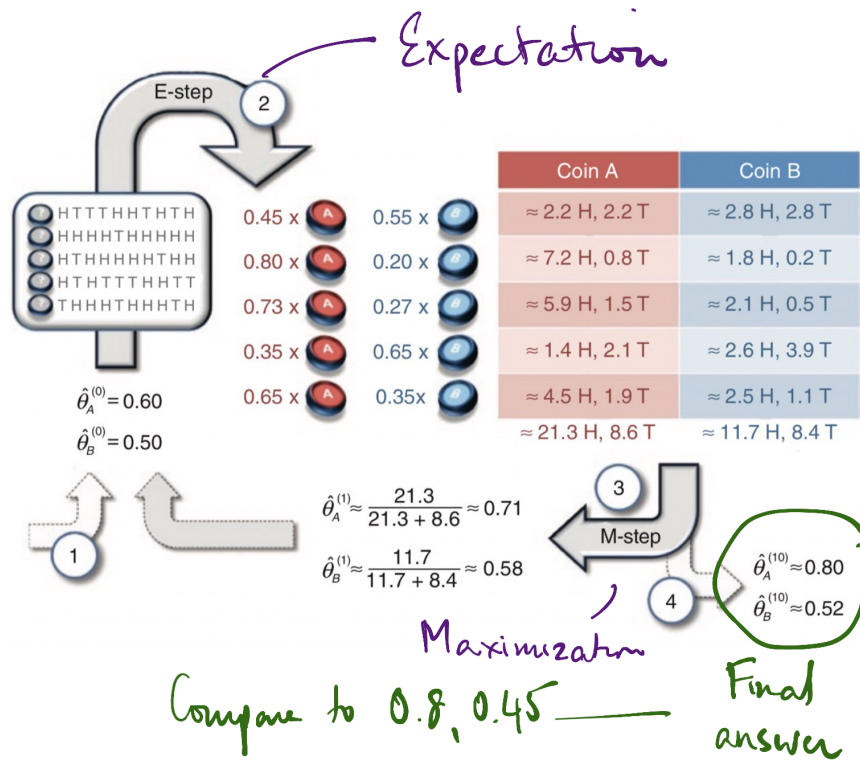
$\frac{P_R}{P_R+P_B} \rightarrow 0.7$ $\frac{P_B}{P_R+P_B} \rightarrow 0.3$

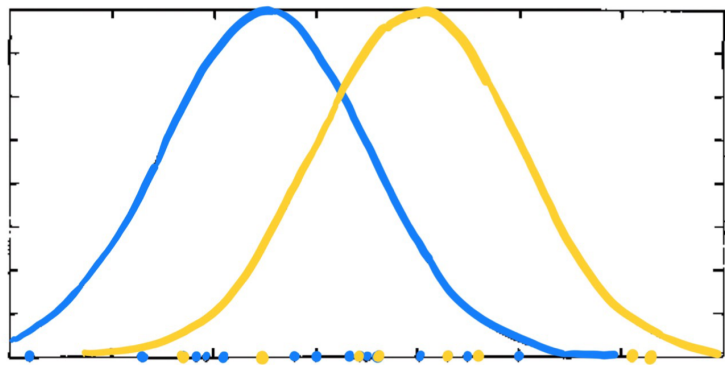
Split the outcome in this ratio

3H, 7T — Red get $0.7 \times 3H, 0.7 \times 7T$
 — Blue get $0.3 \times 3H, 0.3 \times 7T$



recompute θ_A, θ_B





Can prove that EM converges if underlying function is convex

Otherwise - local maximum

Choose μ_1, μ_2

For each x_i

$$\frac{P(x_i|\mu_1)}{\sum_{j=1}^2 P(x_i|\mu_j)}, \frac{P(x_i|\mu_2)}{\sum_{j=1}^2 P(x_i|\mu_j)}$$

For each x_i -

$$w_i^1 = \frac{P(x_i|\mu_1)}{\sum}$$

$$w_i^2 = 1 - w_i^1 = \frac{P(x_i|\mu_2)}{\sum}$$

Sample mean $\frac{1}{n} \sum x_i$

Weighted sample mean

$$\sum_{i=1}^n \frac{w_i^1 x_i}{w_i^1} = \mu_1$$

$$\sum_{i=1}^n \frac{w_i^2 x_i}{w_i^2} = \mu_2$$

Back to topic modelling

What happens if we don't know the topics?

- I like to eat broccoli and bananas.
- I ate a banana and spinach smoothie for breakfast.
- Chinchillas and kittens are cute.
- My sister adopted a kitten yesterday.
- Look at this cute hamster munching on a piece of broccoli.

Told that there are 2 topics $\alpha T_1 + (1-\alpha) T_2$

K topics : $p_1, p_2 \dots p_k$, $\sum p_i = 1$

For each word w_i & topic t_j $P(w_i | t_j)$

Latent Dirichlet Analysis = LDA

Step 1 For each word in the overall set of docs, randomly assign a topic

Step 2 Assign a mixture of topics to each document

For each word w_i overall

For each topic t_j

Estimate $P(w_i | t_j)$ globally

Reassign a topic to w_i according to these probabilities

Re-estimate topic allocations for each doc

Iterate till convergence

Semi-Supervised Learning

Biggest bottleneck of Supervised learning is providing labelled training data