

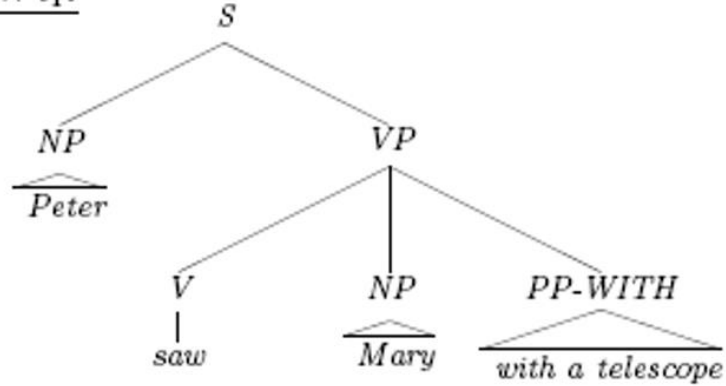
Statistical Parsing: Part 2

An Example

①

$100 \times t_1:$

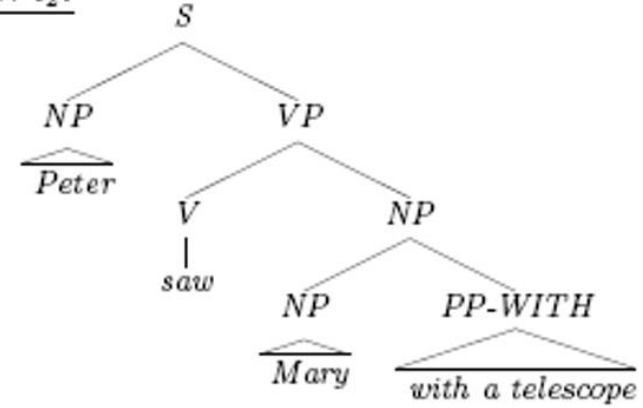
✓



②

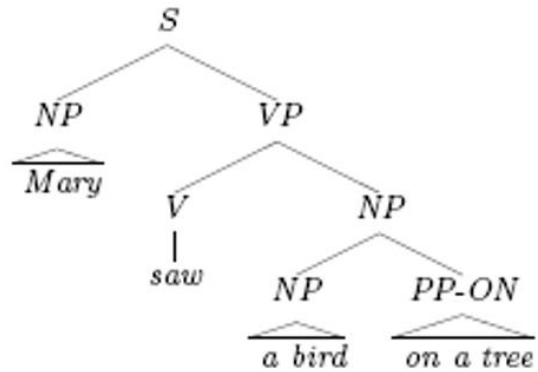
$5 \times t_2:$

✓



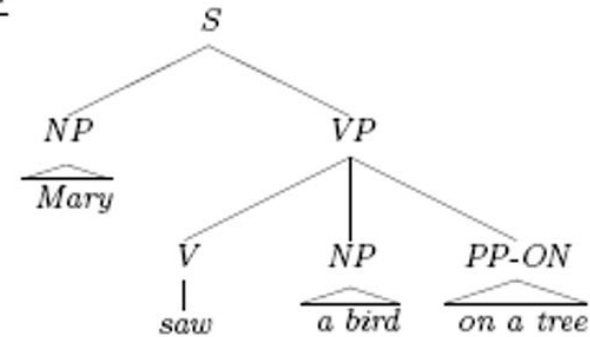
③

$100 \times t_3:$



④

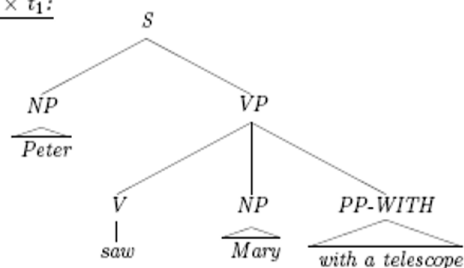
$5 \times t_4:$



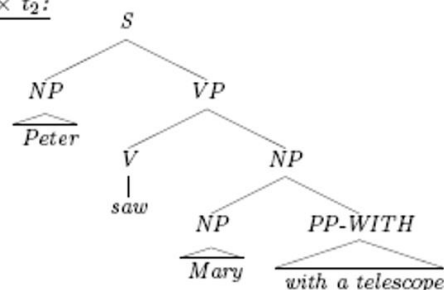
[Refer Workbook](#)

Parameter Estimation: Maximum Likelihood Estimates

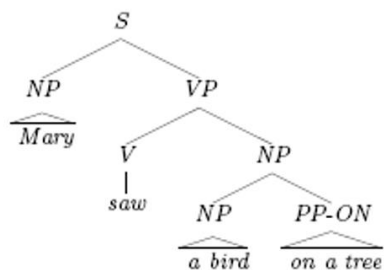
$100 \times t_1:$



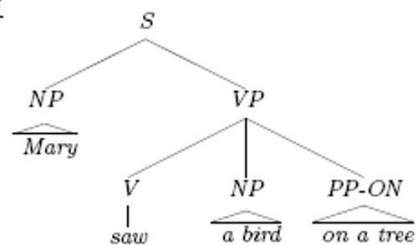
$5 \times t_2:$



$100 \times t_3:$



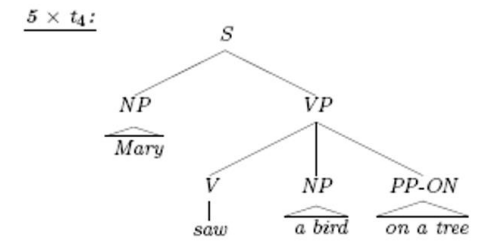
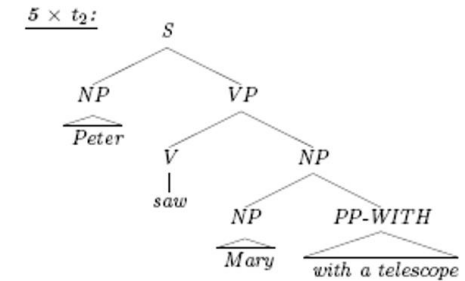
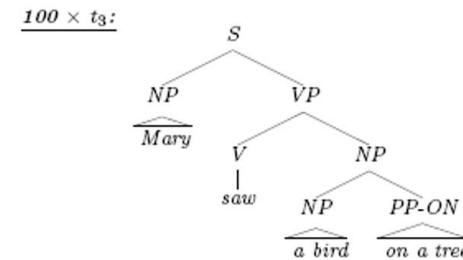
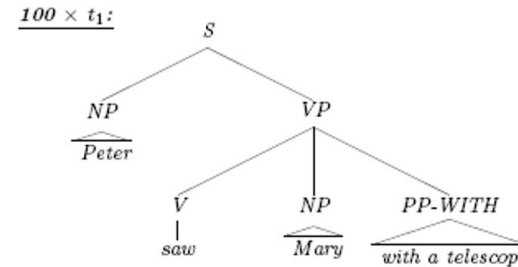
$5 \times t_4:$



CFG rule	Rule frequency	Rule probability
$S \rightarrow NP VP$	$100 + 5 + 100 + 5$	$\frac{210}{210} = 1.000$
$VP \rightarrow V NP PP-WITH$	100	$\frac{100}{210} \approx 0.476$
$VP \rightarrow V NP PP-ON$	5	$\frac{5}{210} \approx 0.024$
$VP \rightarrow V NP$	$5 + 100$	$\frac{105}{210} = 0.500$
$NP \rightarrow Peter$	$100 + 5$	$\frac{105}{525} = 0.200$
$NP \rightarrow Mary$	$100 + 5 + 100 + 5$	$\frac{210}{525} = 0.400$
$NP \rightarrow a \text{ bird}$	$100 + 5$	$\frac{105}{525} = 0.200$
$NP \rightarrow NP PP-WITH$	5	$\frac{5}{525} \approx 0.010$
$NP \rightarrow NP PP-ON$	100	$\frac{100}{525} \approx 0.190$
$PP-WITH \rightarrow \text{with a telescope}$	$100 + 5$	$\frac{105}{105} = 1.000$
$PP-ON \rightarrow \text{on a tree}$	$100 + 5$	$\frac{105}{105} = 1.000$
$V \rightarrow \text{saw}$	$100 + 5 + 100 + 5$	$\frac{210}{210} = 1.000$

Using the PCFG for disambiguation

CFG rule	Rule frequency	Rule probability
$S \rightarrow NP VP$	$100 + 5 + 100 + 5$	$\frac{210}{210} = 1.000$
$VP \rightarrow V NP PP-WITH$	100	$\frac{100}{210} \approx 0.476$
$VP \rightarrow V NP PP-ON$	5	$\frac{5}{210} \approx 0.024$
$VP \rightarrow V NP$	$5 + 100$	$\frac{105}{210} = 0.500$
$NP \rightarrow \text{Peter}$	$100 + 5$	$\frac{105}{525} = 0.200$
$NP \rightarrow \text{Mary}$	$100 + 5 + 100 + 5$	$\frac{210}{525} = 0.400$
$NP \rightarrow \text{a bird}$	$100 + 5$	$\frac{105}{525} = 0.200$
$NP \rightarrow NP PP-WITH$	5	$\frac{5}{525} \approx 0.010$
$NP \rightarrow NP PP-ON$	100	$\frac{100}{525} \approx 0.190$
$PP-WITH \rightarrow \text{with a telescope}$	$100 + 5$	$\frac{105}{105} = 1.000$
$PP-ON \rightarrow \text{on a tree}$	$100 + 5$	$\frac{105}{105} = 1.000$
$V \rightarrow \text{saw}$	$100 + 5 + 100 + 5$	$\frac{210}{210} = 1.000$



"Peter saw Mary with a telescope"

$$p(VP \rightarrow V NP PP-WITH) > p(VP \rightarrow V NP) \cdot p(NP \rightarrow NP PP-WITH)$$

$$0.476 > 0.500 \cdot 0.010$$

"Mary saw a bird on a tree"

$$p(VP \rightarrow V NP PP-ON) < p(VP \rightarrow V NP) \cdot p(NP \rightarrow NP PP-ON)$$

$$0.024 < 0.500 \cdot 0.190$$

Try this : Mary saw a bird on a tree with a telescope

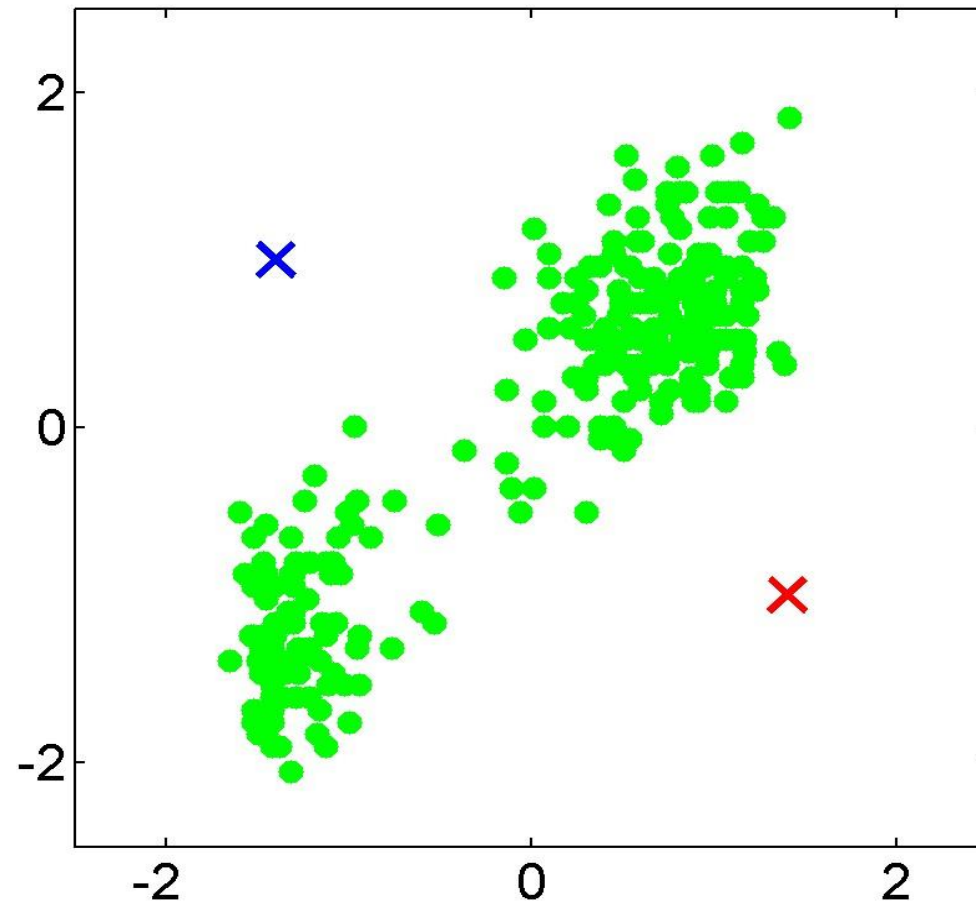
Incomplete data problem in PCFG parameter estimation

$f(y)$	y	
5	y_1	$y_1 = \text{"Mary saw a bird on a tree"}$
10	y_2	$y_2 = \text{"a bird on a tree saw a worm"}$

$S \longrightarrow NP VP$
 $VP \longrightarrow V NP$
 $VP \longrightarrow V NP PP$
 $NP \longrightarrow NP PP$
 $NP \longrightarrow \text{Mary}$
 $NP \longrightarrow \text{a bird}$
 $NP \longrightarrow \text{a worm}$
 $PP \longrightarrow \text{on a tree}$
 $V \longrightarrow \text{saw}$

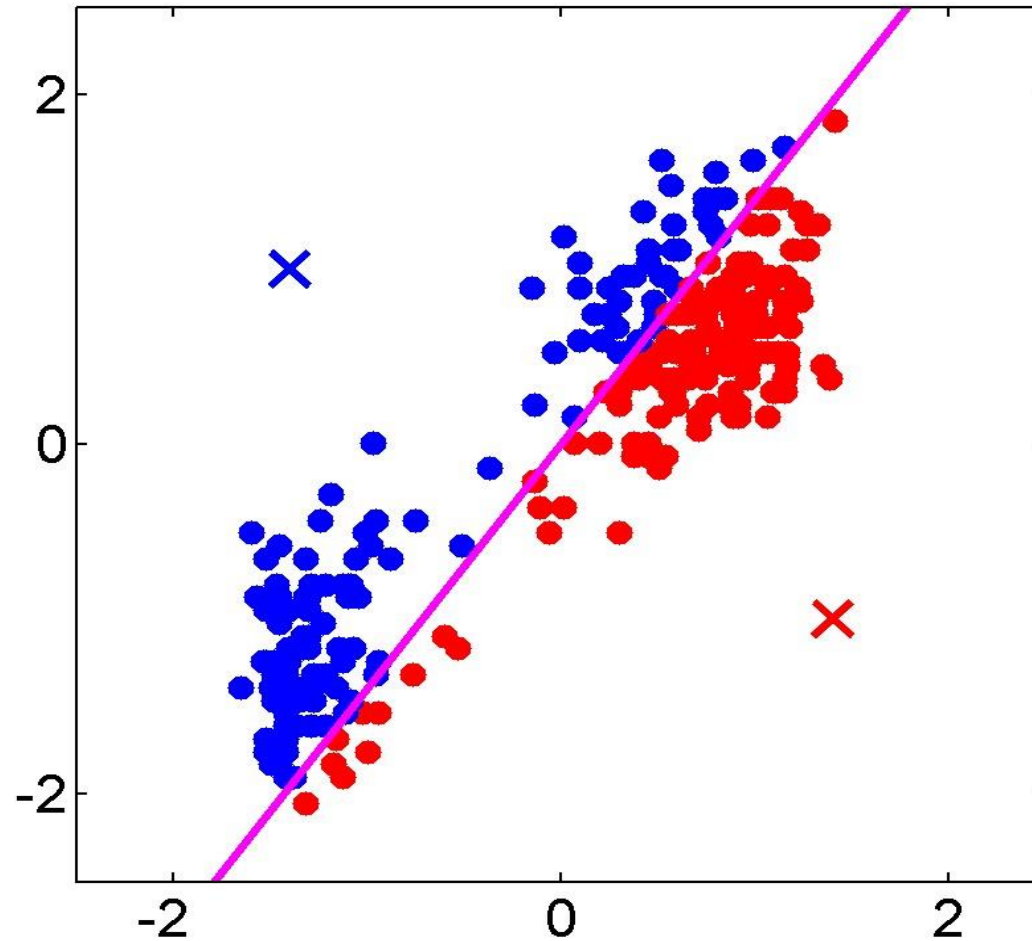
[Refer Workbook](#)

K-means clustering: M_0



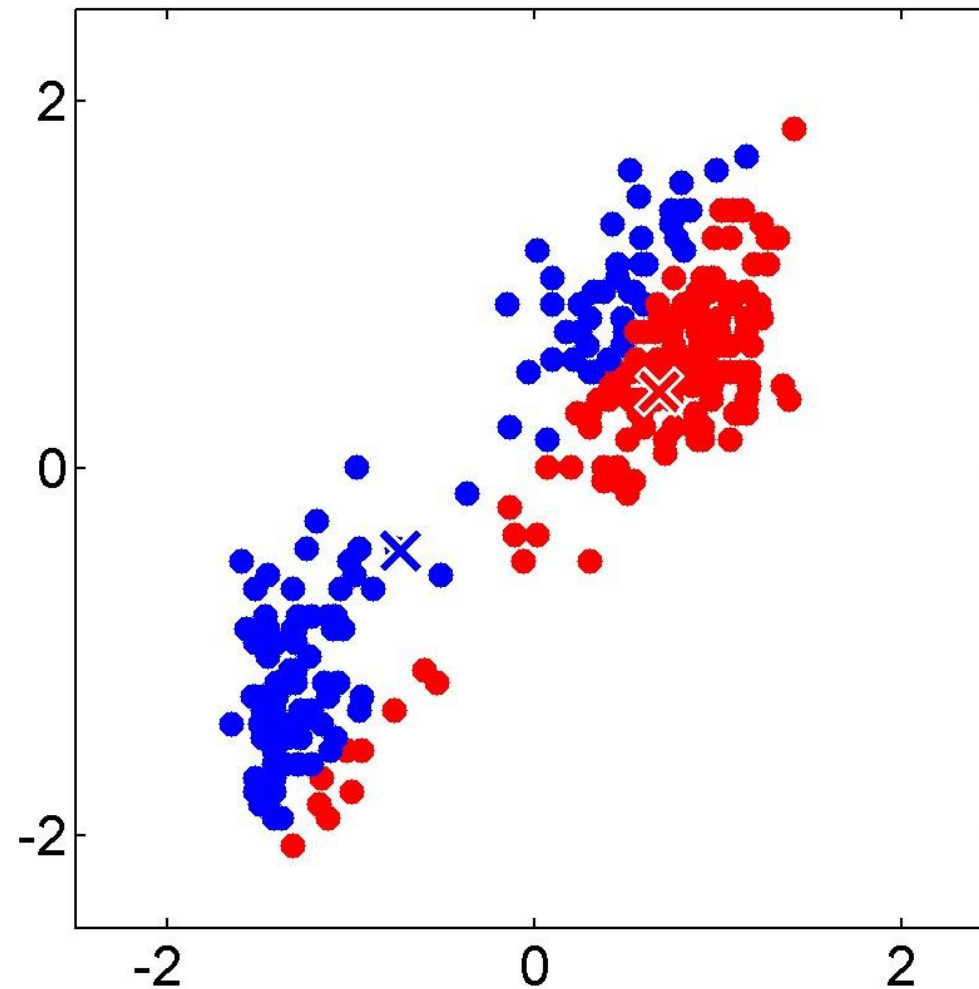
Ack: Chris Bishop

K-means clustering: E_1



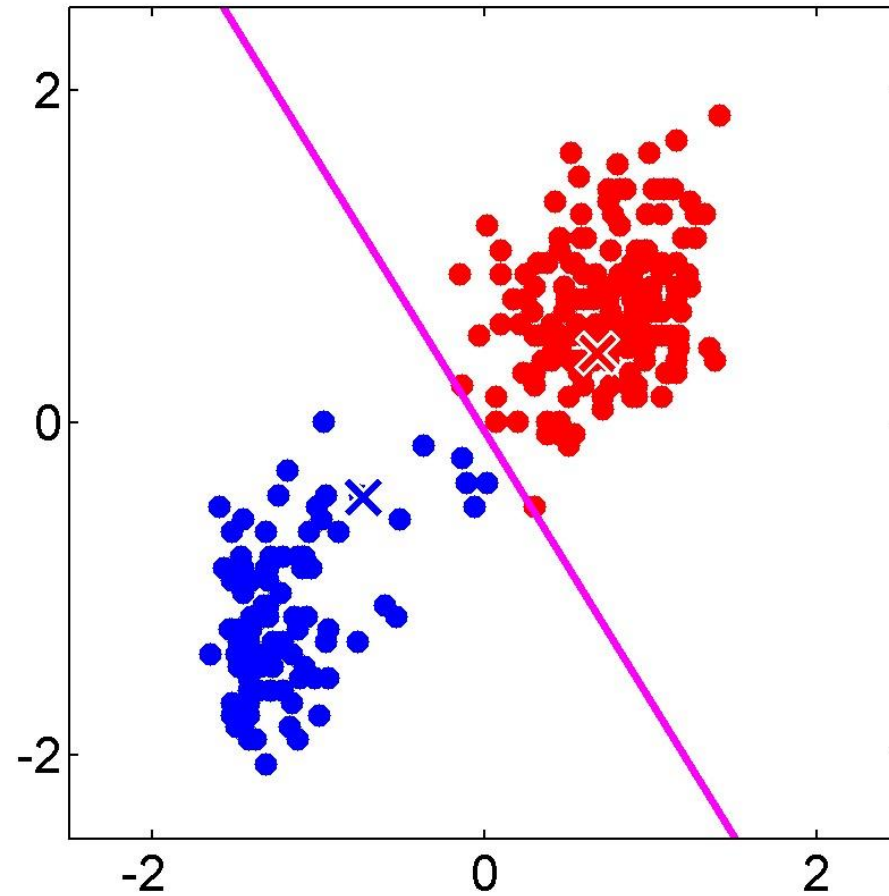
Ack: Chris Bishop

K-means clustering: M_1



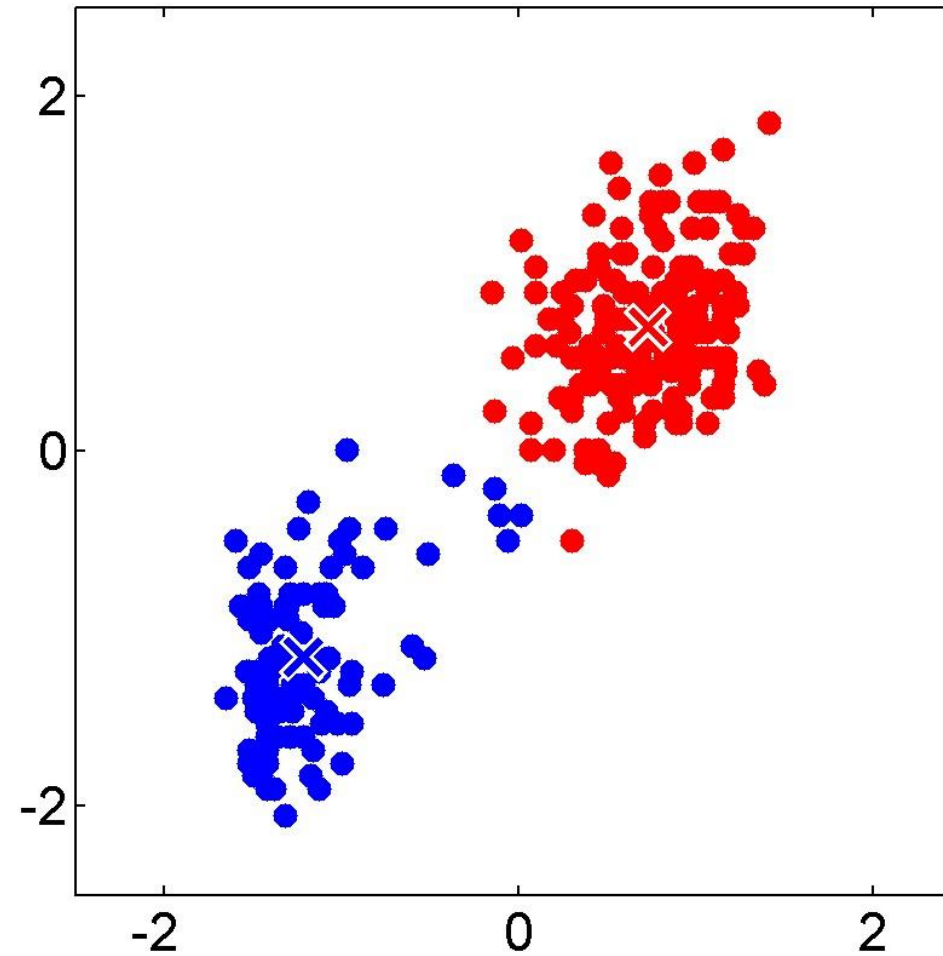
Ack: Chris Bishop

K-means clustering: E_2



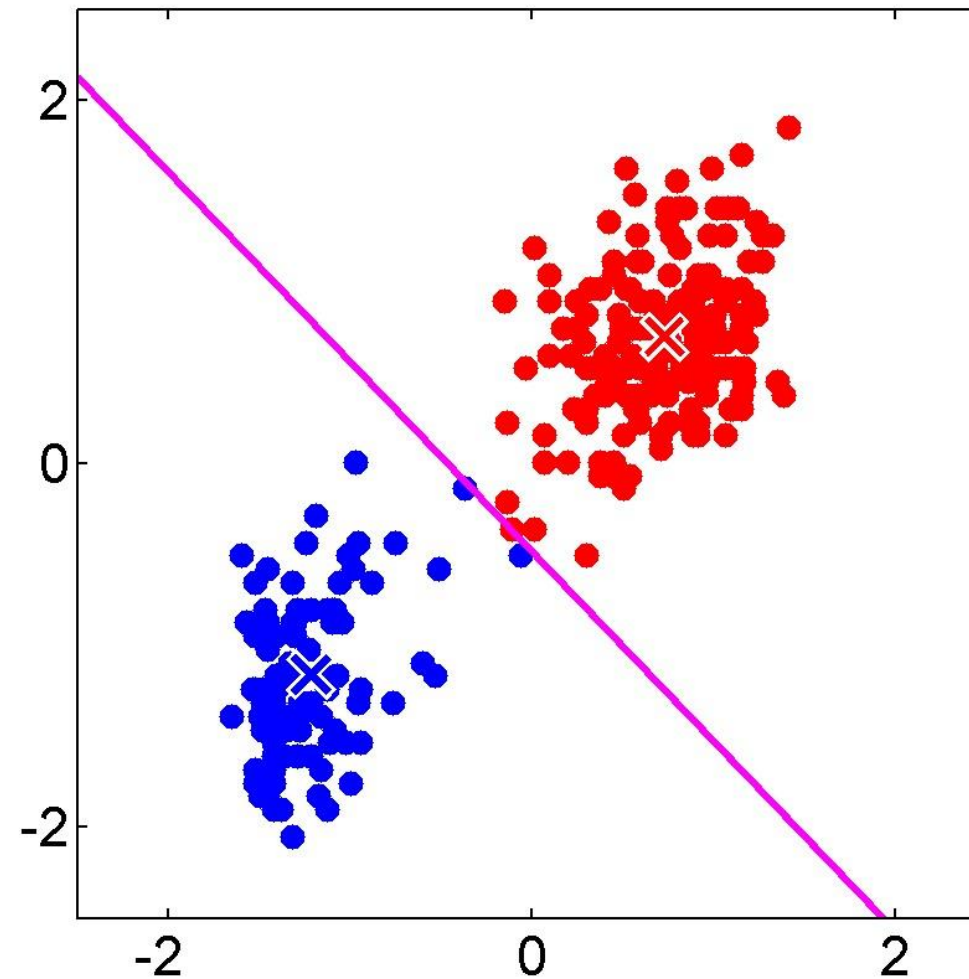
Ack: Chris Bishop

K-means clustering: M_2



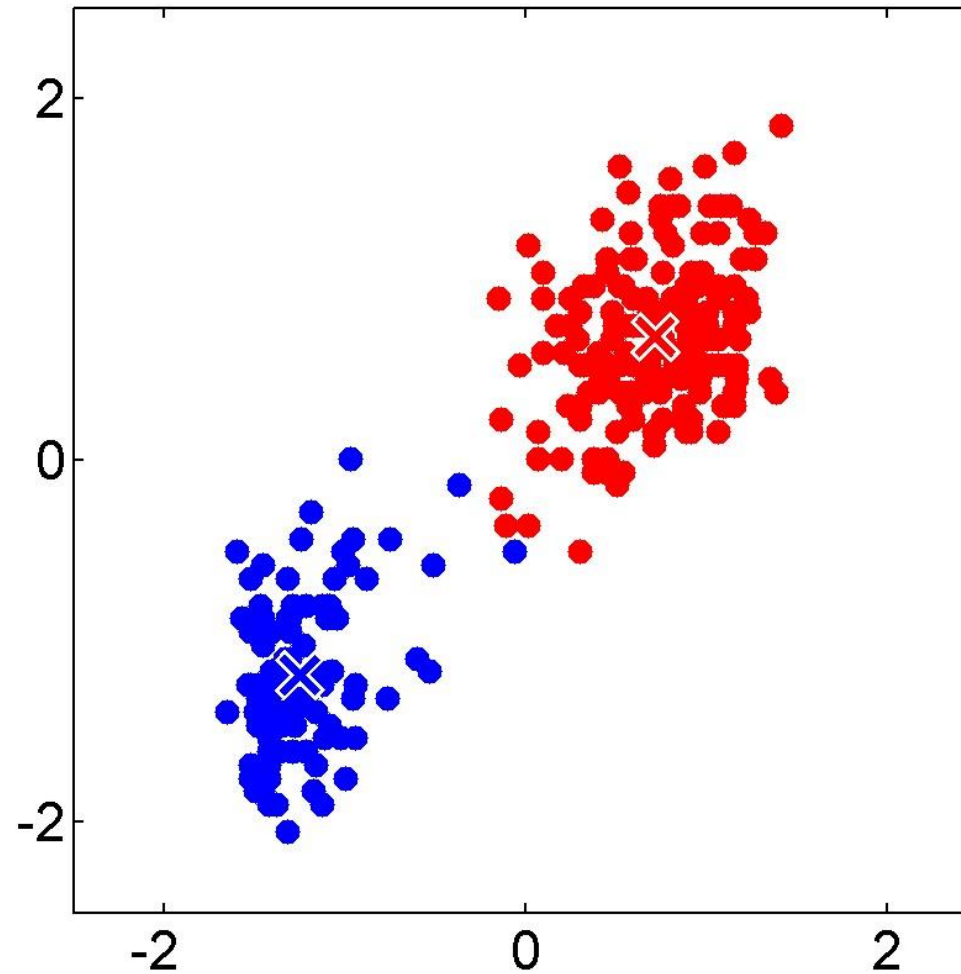
Ack: Chris Bishop

K-means clustering: E_3



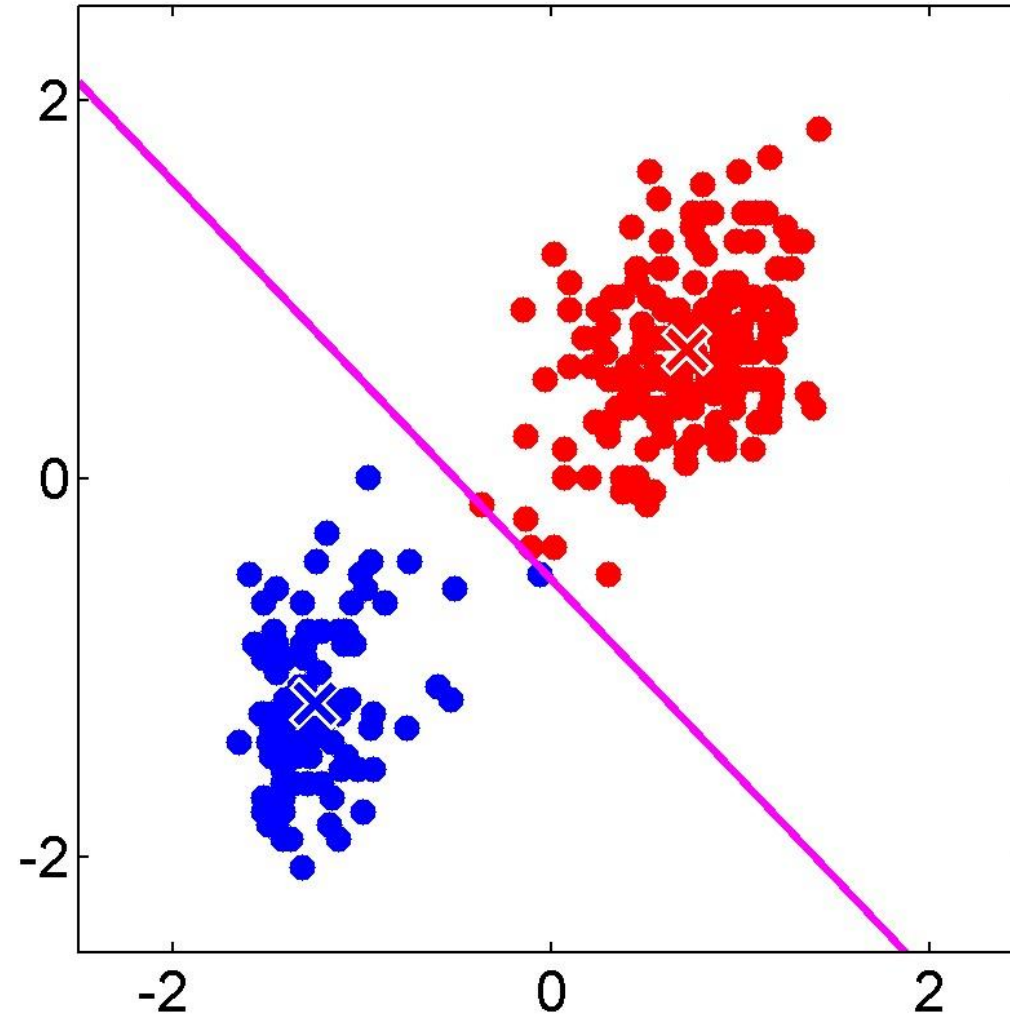
Ack: Chris Bishop

K-means clustering: M_3



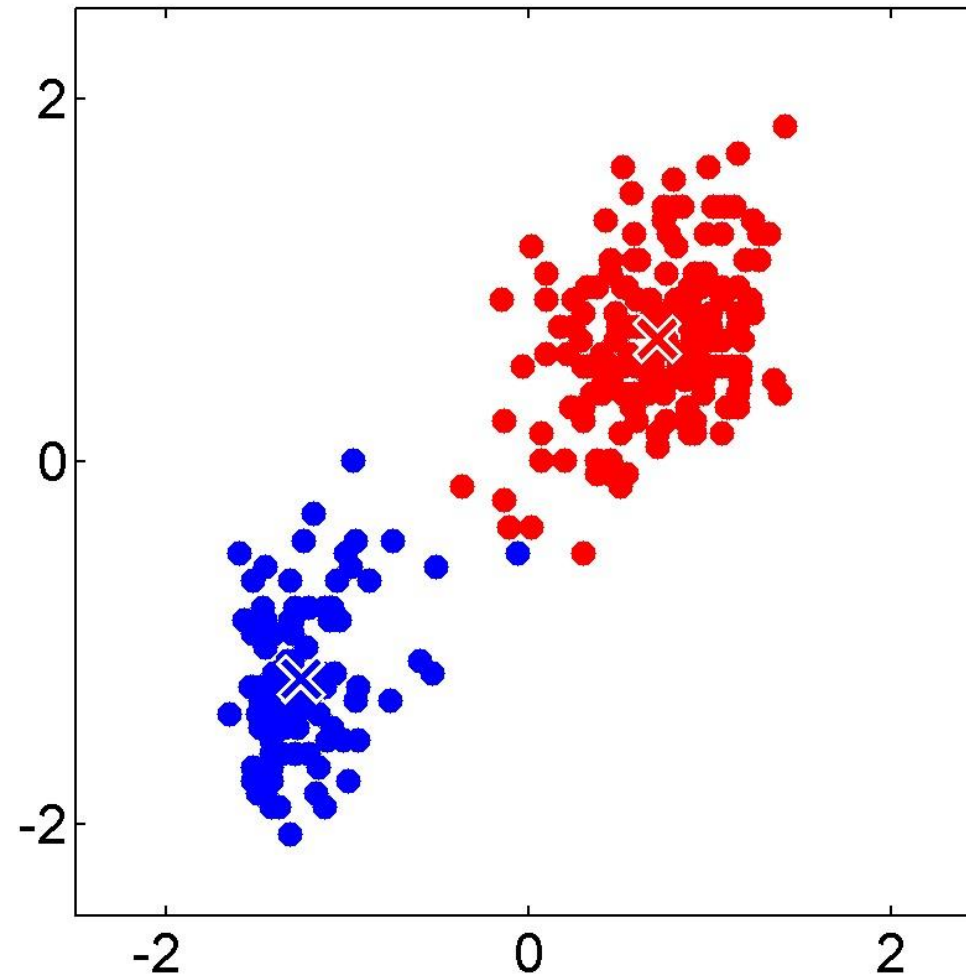
Ack: Chris Bishop

K-means clustering: E_4



Ack: Chris Bishop

K-means clustering: M_4



Ack: Chris Bishop

Responsibilities

- *Responsibilities* assign data points to clusters $r_{nk} \in \{0, 1\}$

such that $\sum_k r_{nk} = 1$

Example: 5 data points and 3 clusters

$$(r_{nk}) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

Ack: Chris Bishop

K-means Cost Function

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

data

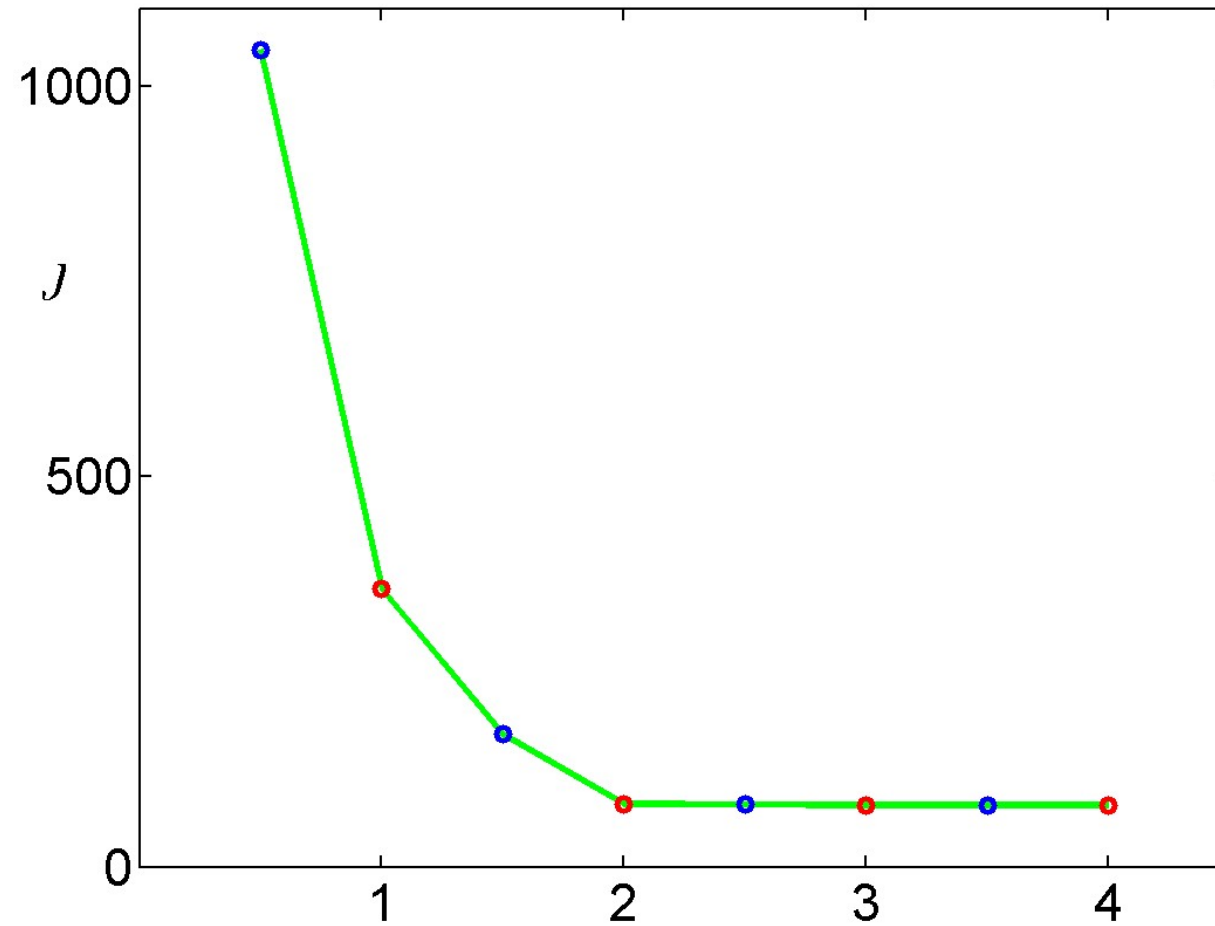
responsibilities

prototypes

The diagram shows the K-means cost function $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$. Three blue arrows point from text labels to parts of the equation: 'data' points to \mathbf{x}_n , 'responsibilities' points to r_{nk} , and 'prototypes' points to $\boldsymbol{\mu}_k$.

Ack: Chris Bishop

Minimizing the Cost Function



Ack: Chris Bishop

Minimizing the Cost Function

- E-step: minimize J w.r.t. r_{nk}
 - assigns each data point to nearest prototype
- M-step: minimize J w.r.t μ_k
 - gives

$$\mu_k = \frac{\sum_n r_{kn} \mathbf{x}_n}{\sum_n r_{kn}}$$

- each prototype set to the mean of points in that cluster
- Convergence guaranteed since there is a finite number of possible settings for the responsibilities

Ack: Chris Bishop

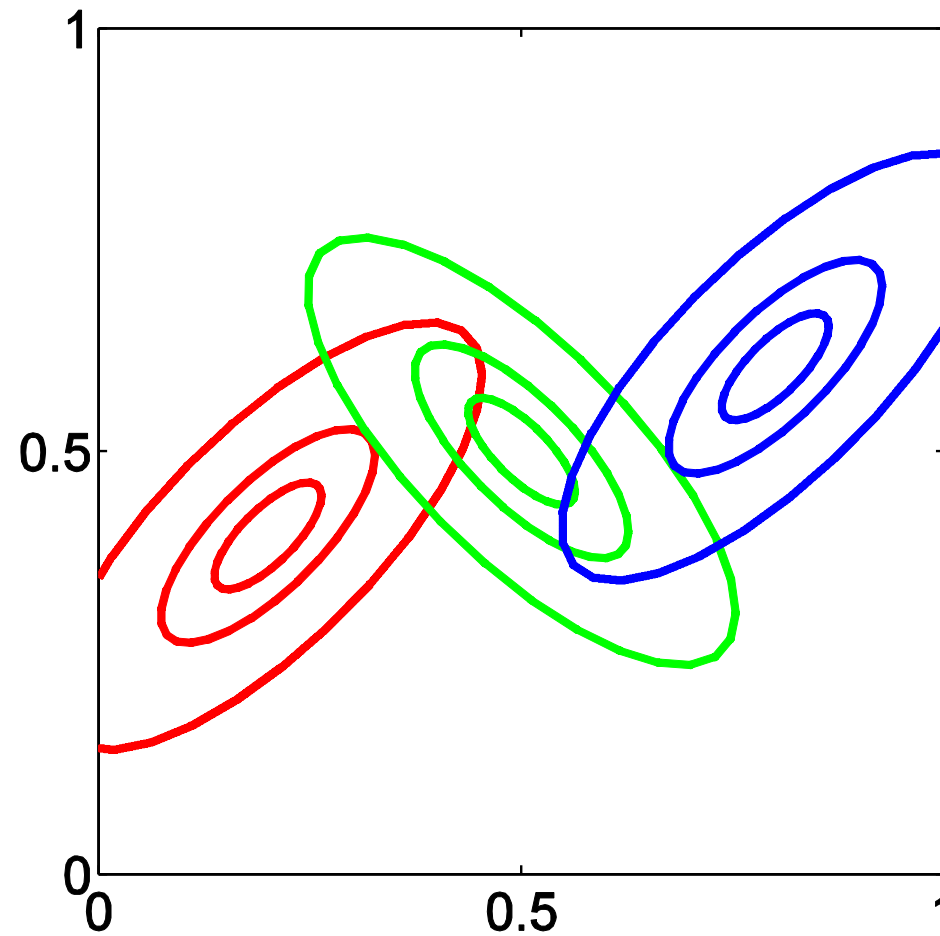
Observation

Limitation: Hard assignments of data points to clusters – small shift of a data point can flip it to a different cluster

Solution: replace 'hard' clustering of K-means with 'soft' probabilistic assignments
Represent the probability distribution of the data as a *Gaussian Mixture Model*

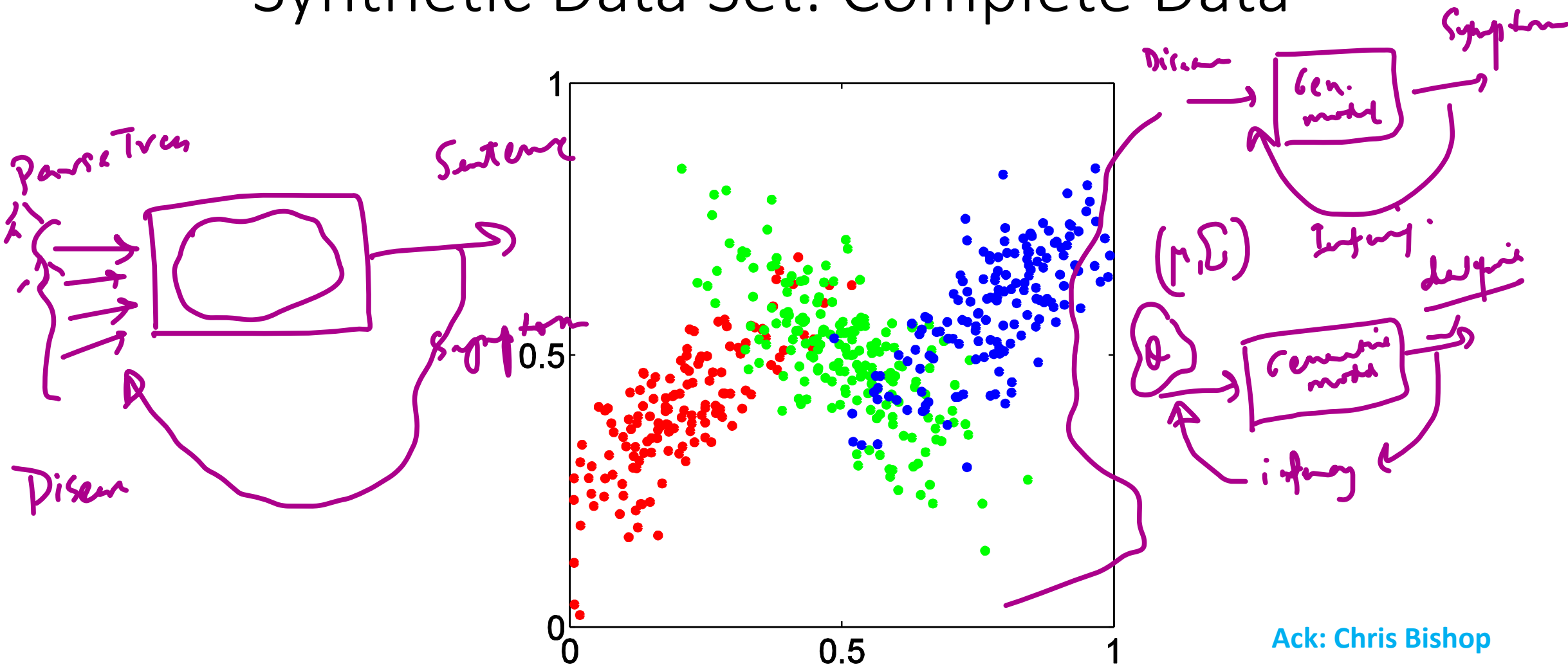
Ack: Chris Bishop

Gaussian Mixture Model



Ack: Chris Bishop

Synthetic Data Set: Complete Data



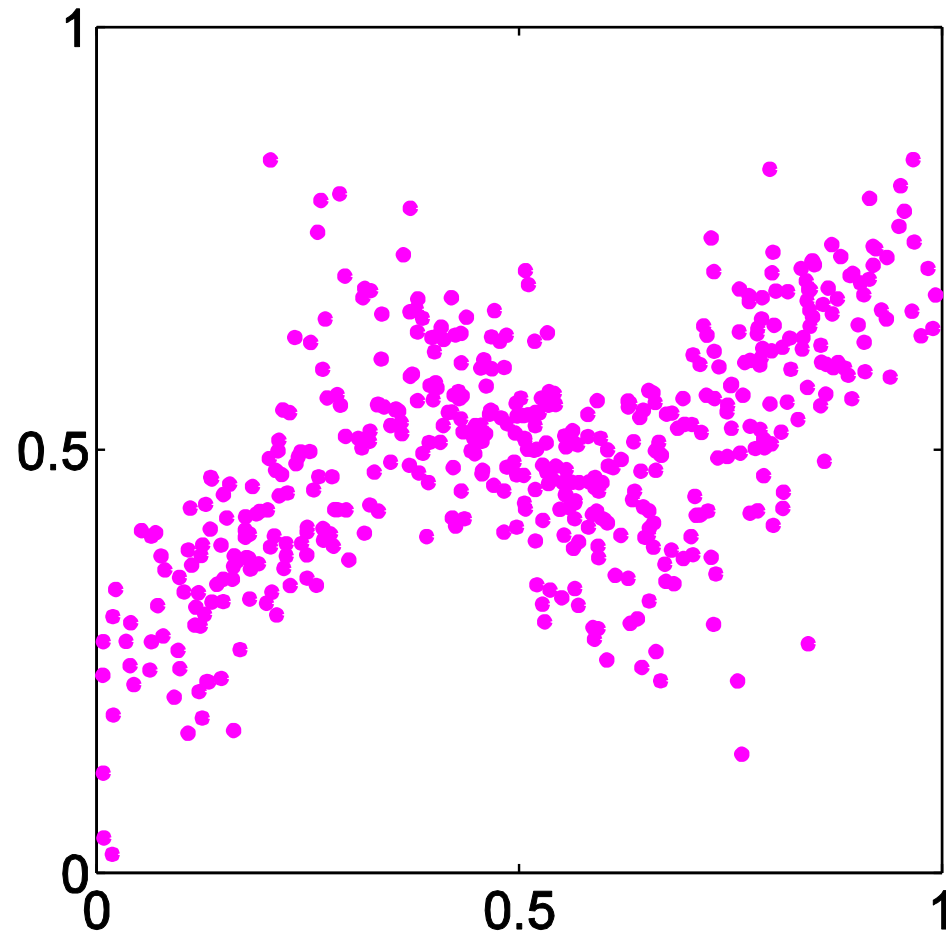
Ack: Chris Bishop

Fitting the Gaussian Mixture

- We wish to invert this process – given the data set, find the corresponding parameters:
 - mixing coefficients
 - means
 - covariances
- If we knew which component generated each data point, the maximum likelihood solution would involve fitting each component to the corresponding cluster
- Problem: the data set is unlabelled
- We shall refer to the labels as *latent* (= hidden) variables

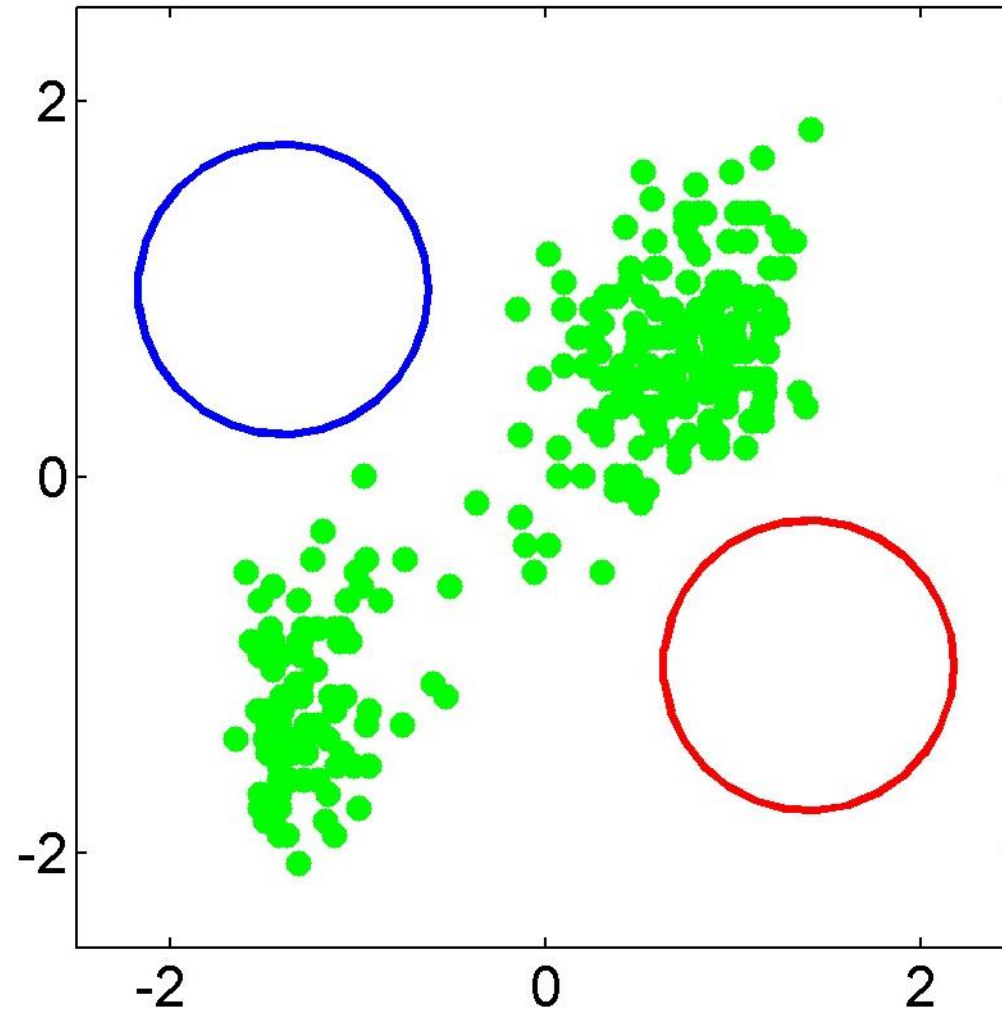
Ack: Chris Bishop

Synthetic Data Set Without Labels



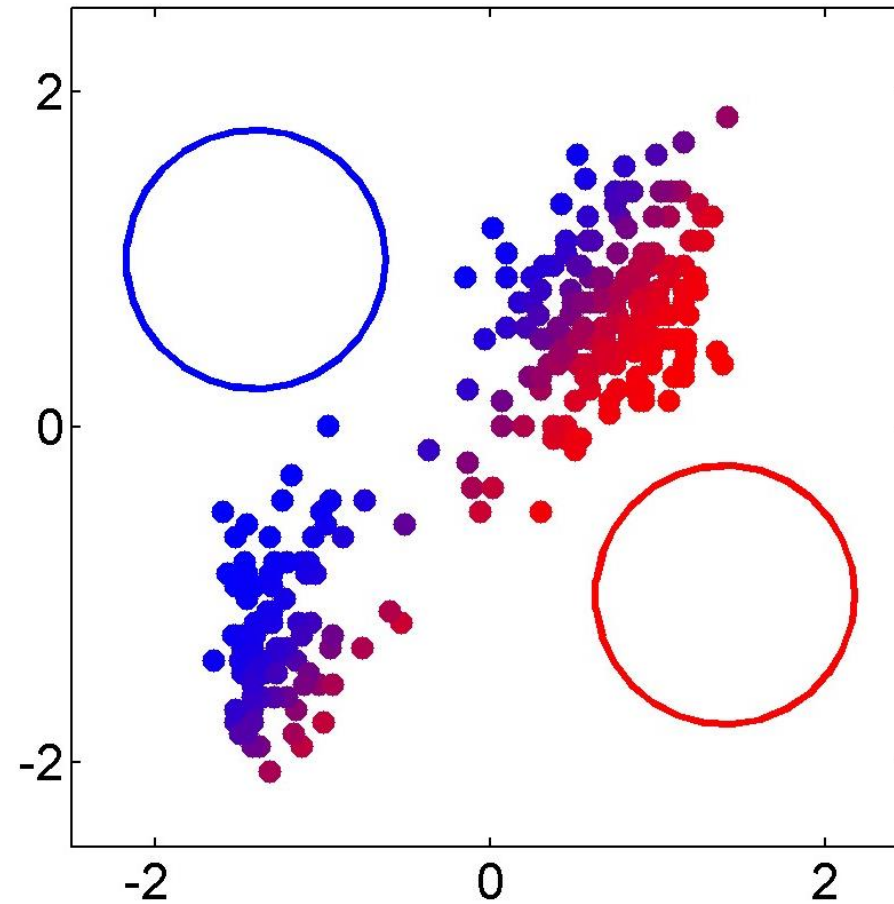
Ack: Chris Bishop

Synthetic Data: Incomplete Data (M_0)



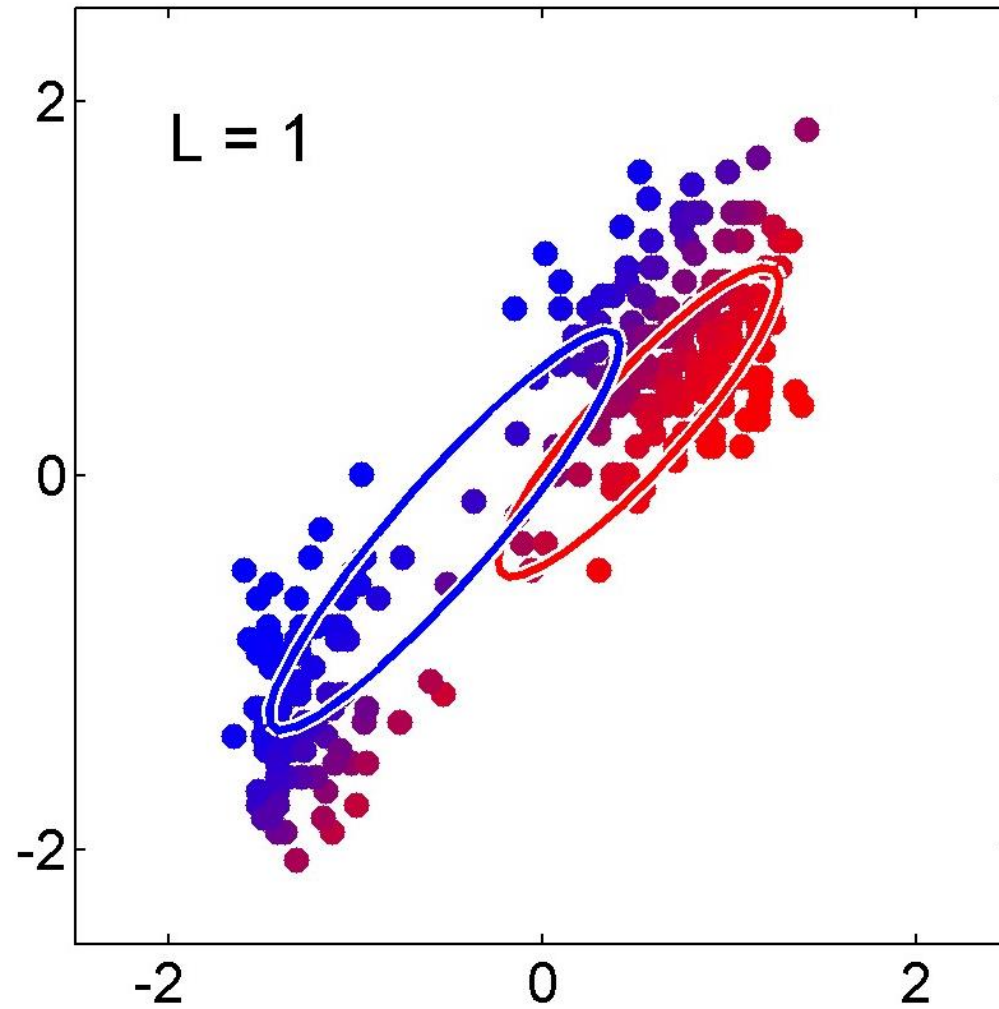
Ack: Chris Bishop

GMM: $M_0 + E_1$



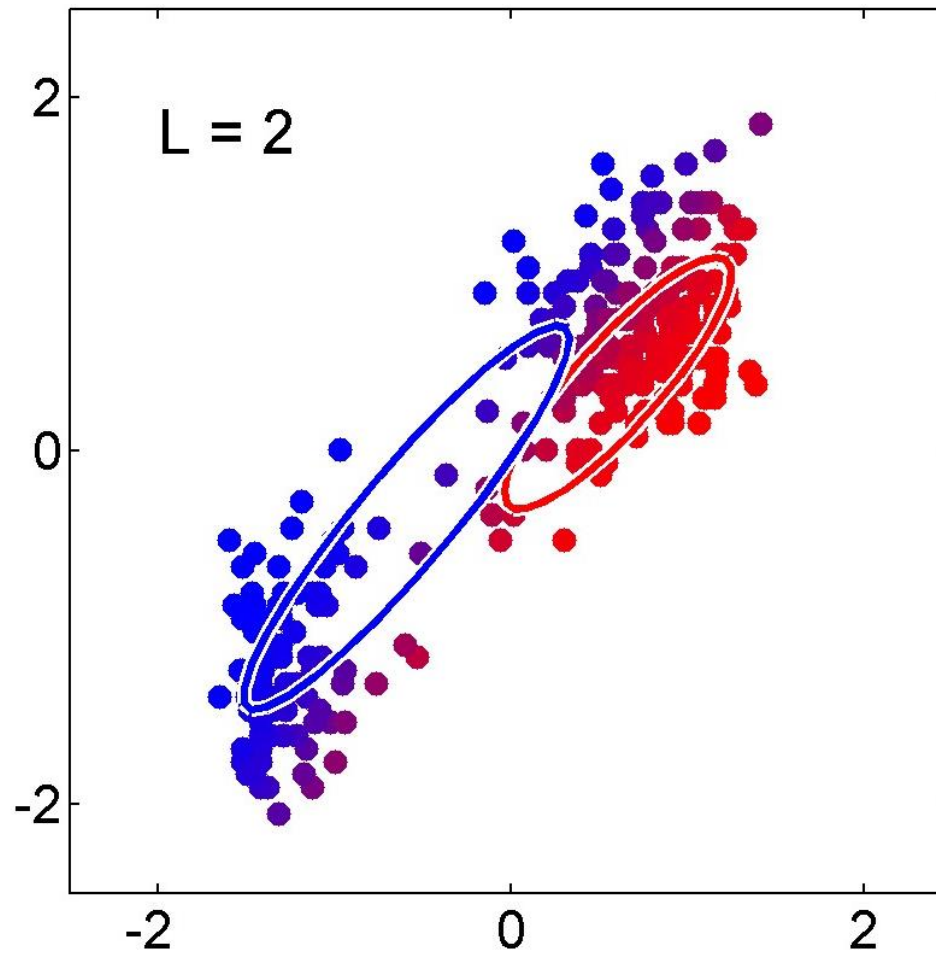
Ack: Chris Bishop

GMM: M_1



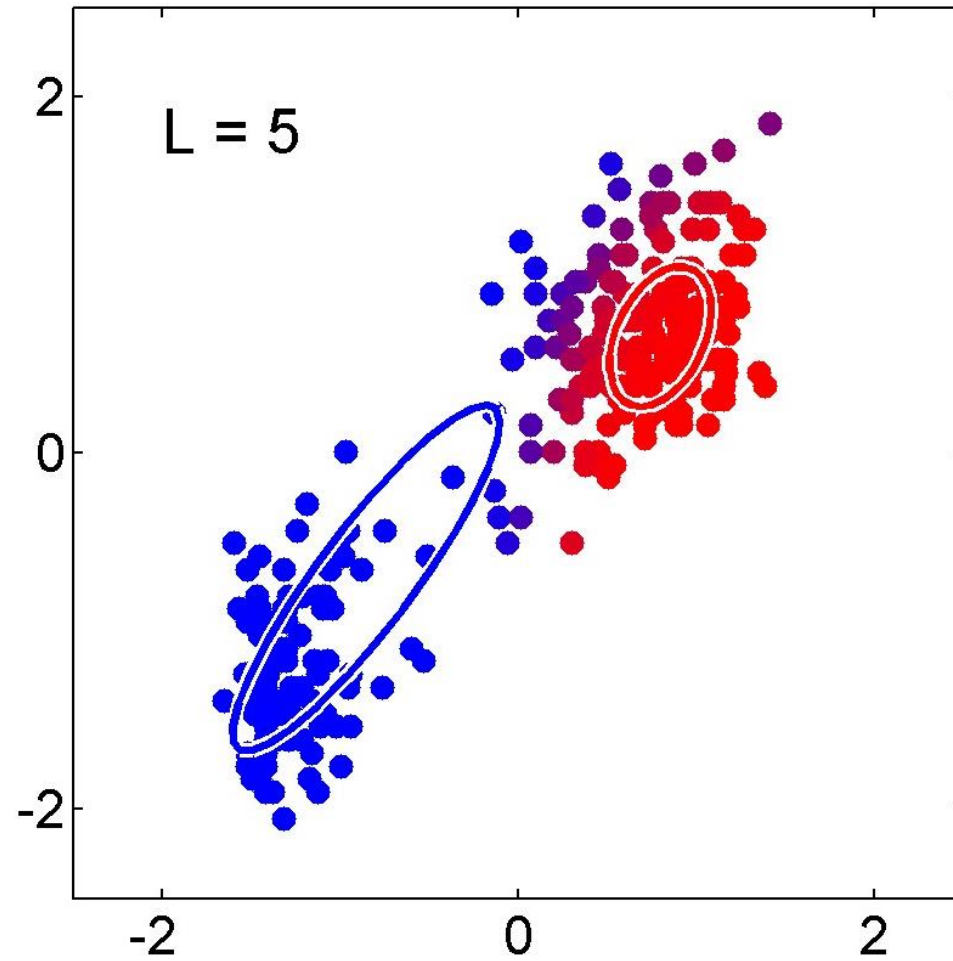
Ack: Chris Bishop

GMM: $M_1 + E_2$



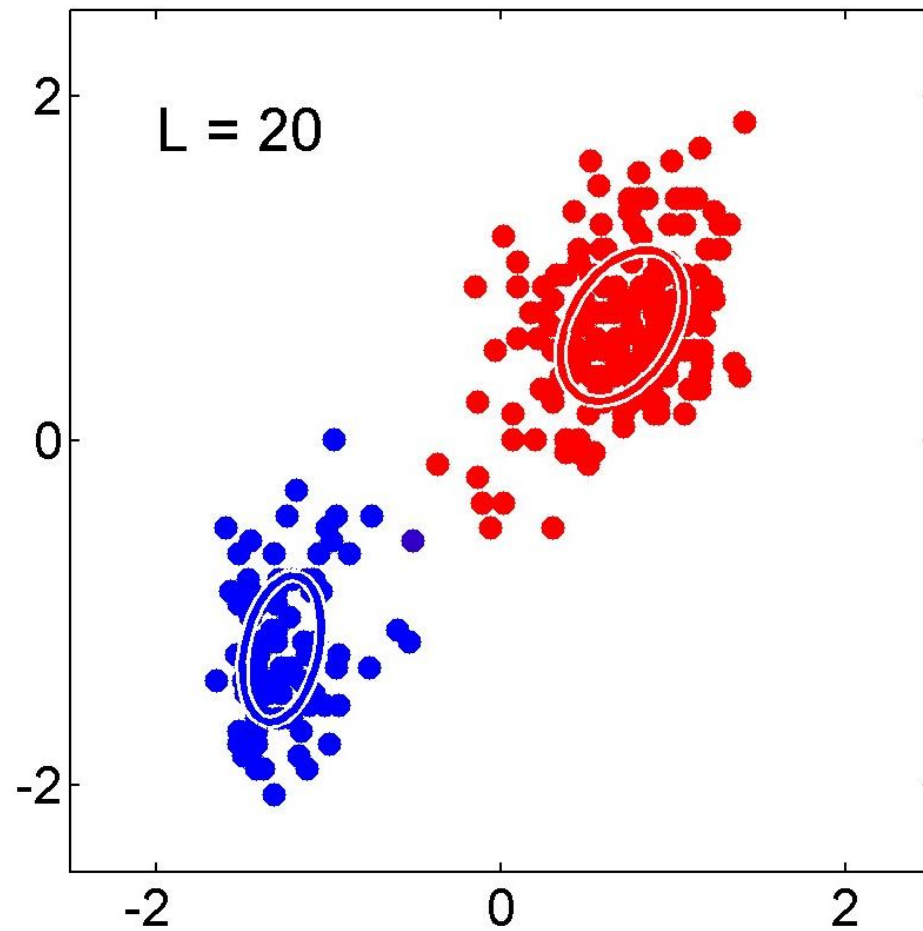
Ack: Chris Bishop

GMM: $M_2 + E_3$



Ack: Chris Bishop

GMM: $M_3 + E_4$



Ack: Chris Bishop

Example 2: Complete data

a Maximum likelihood

	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

5 sets, 10 tosses per set

Coin A	Coin B
	5 H, 5 T
9 H, 1 T	
8 H, 2 T	
	4 H, 6 T
7 H, 3 T	
24 H, 6 T	9 H, 11 T

$$\theta_A^k (1-\theta)^{n-k}$$

$\left(\frac{k}{n}\right)$

$$\hat{\theta}_A = \frac{24}{24 + 6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9 + 11} = 0.45$$

Refer: shared paper

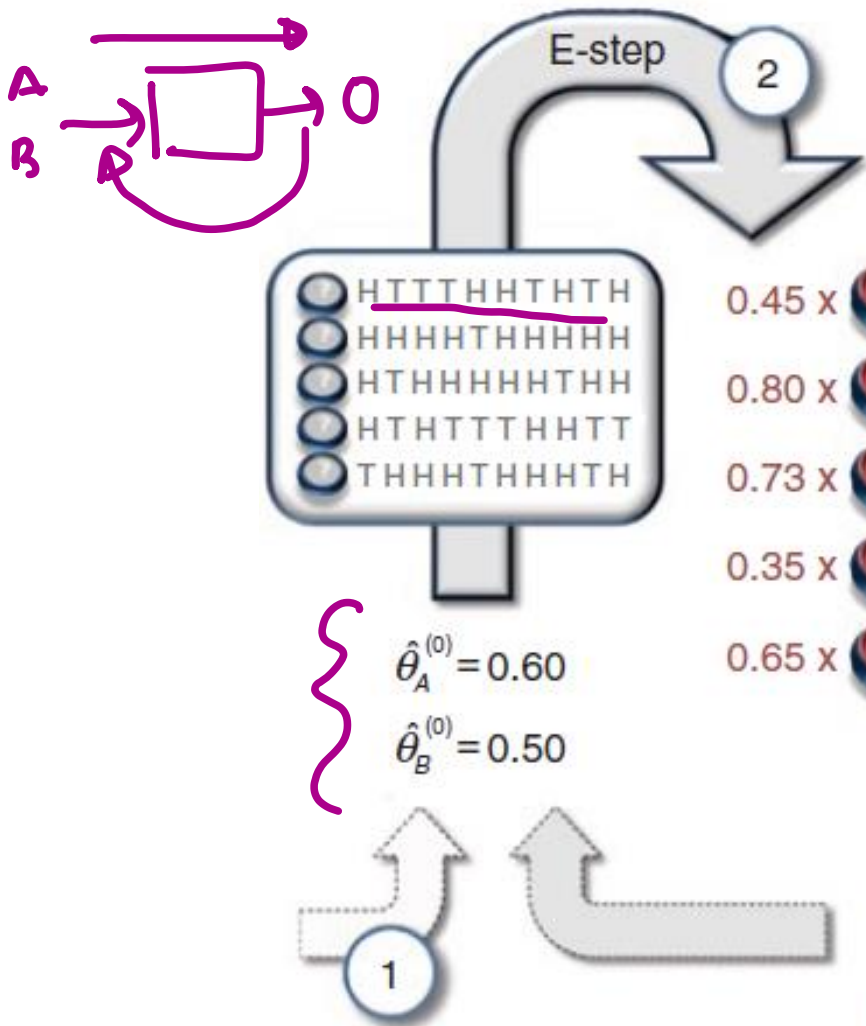
Example 2: Incomplete Data

$$P(A|O) = \frac{P(O|A) \cdot P(A)}{P(O|A) \cdot P(A) + P(O|B) \cdot P(B)}$$

$$= \frac{(0.6)^5 (0.4)^5}{(0.6)^5 (0.4)^5 + (0.5)^5}$$

$$= 0.45$$

Max^m Likelihood

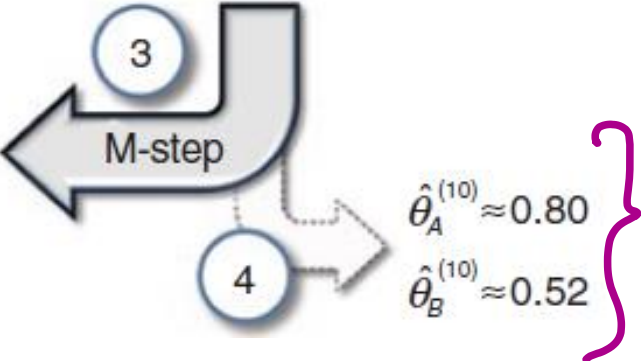


0.45 x	A	0.55 x	B
0.80 x	A	0.20 x	B
0.73 x	A	0.27 x	B
0.35 x	A	0.65 x	B
0.65 x	A	0.35 x	B

Coin A	Coin B
≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
≈ 1.4 H, 2.1 T	≈ 2.6 H, 3.9 T
≈ 4.5 H, 1.9 T	≈ 2.5 H, 1.1 T
≈ 21.3 H, 8.6 T	≈ 11.7 H, 8.4 T

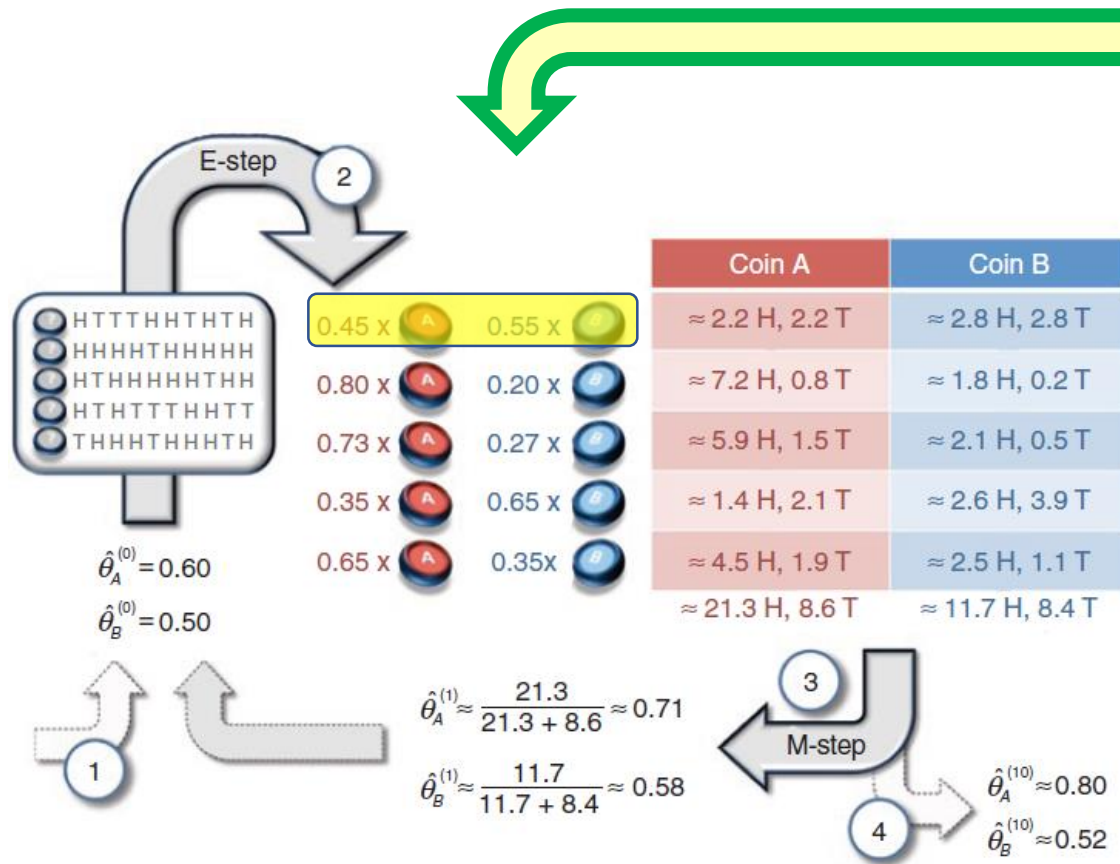
$$\hat{\theta}_A^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{\theta}_B^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$



Refer: shared paper

Example 2: Doubt posted in feedback



$O = \langle H T T T H H T H T H \rangle$

Posterior calculation. 5 heads, 5 tails, used generative model.

$$P(A|O) = \frac{P(O|A) \cdot P(A)}{P(O)} = \frac{P(O|A) \cdot P(A)}{P(O|A) \cdot P(A) + P(O|B) \cdot P(B)}$$

uniform prior.

$$= \frac{(0.6)^5 (0.4)^5 \cdot 0.5}{(0.6)^5 (0.4)^5 \cdot 0.5 + (0.5)^5 (0.5)^5 \cdot 0.5}$$

$$= 0.45$$

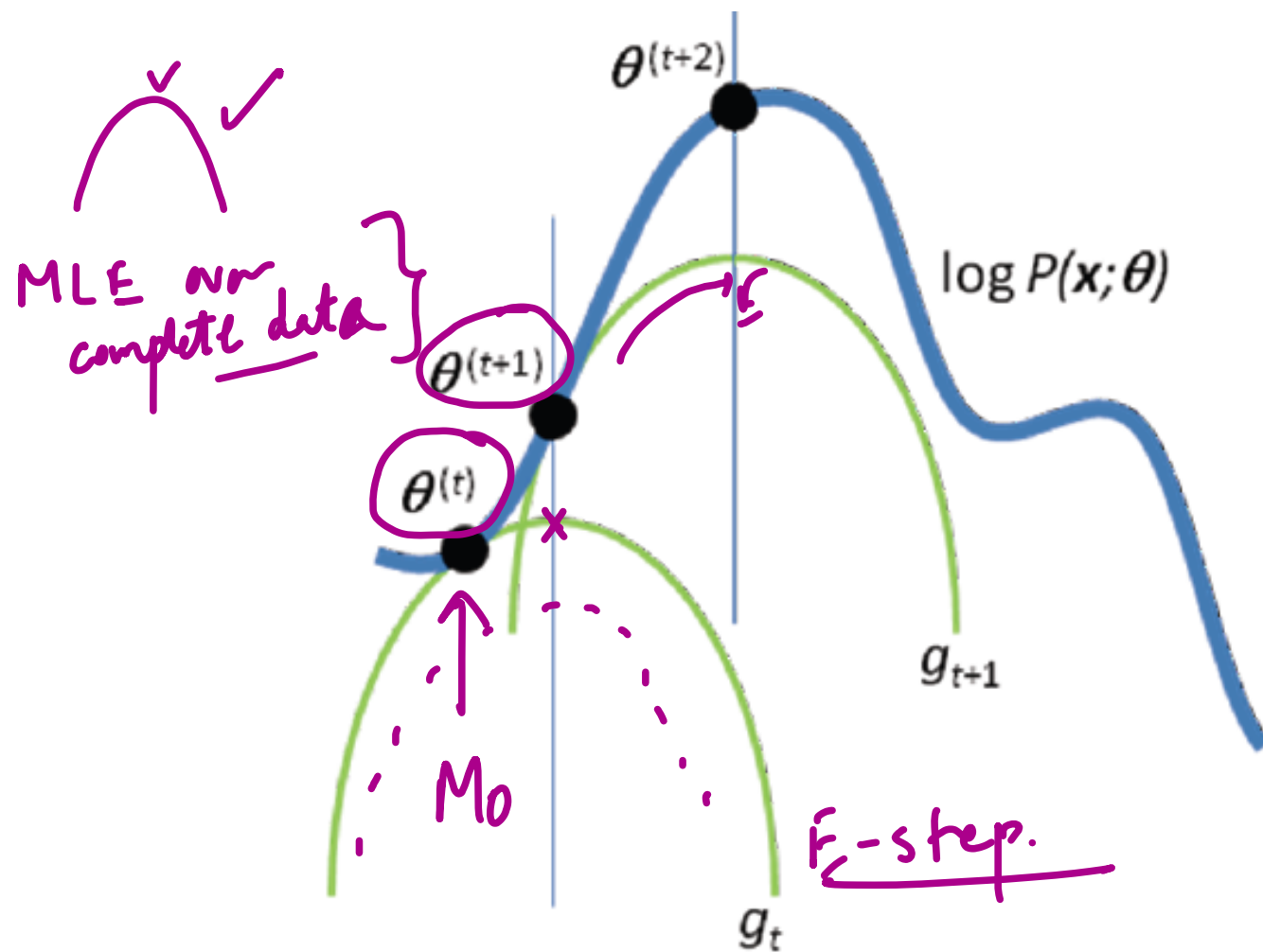
Similarly $P(B|O) = \frac{P(O|B) \cdot P(B)}{P(O|A) \cdot P(A) + P(O|B) \cdot P(B)} = 0.55$

Refer: shared paper

Goal of E step is to estimate posteriors like $P(A|O)$.

Estimating $P(O|A)$ is the preparation for E step. This involves using the generative storyline.

EM intuition : lower bound maximization



Refer: shared paper

EM intuition



EM over PCFGs : an example

$f(y)$	y
5	y_1
10	y_2

$y_1 = \text{"Mary saw a bird on a tree"}$

$y_2 = \text{"a bird on a tree saw a worm"}$

M_0 step

$S \rightarrow NP VP \quad (1.00)$

$VP \rightarrow V NP \quad (0.50)$

$VP \rightarrow V NP PP \quad (0.50)$

$NP \rightarrow NP PP \quad (0.25)$

$NP \rightarrow \text{Mary} \quad (0.25)$

$NP \rightarrow \text{a bird} \quad (0.25)$

$NP \rightarrow \text{a worm} \quad (0.25)$

$PP \rightarrow \text{on a tree} \quad (1.00)$

$V \rightarrow \text{saw} \quad (1.00)$

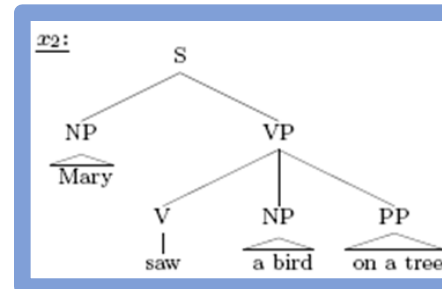
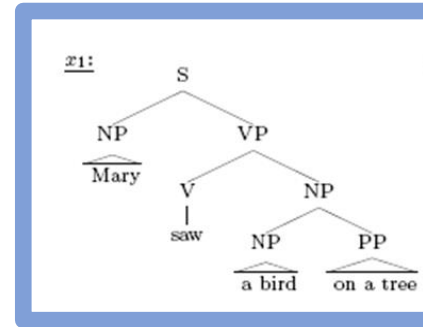
[Refer Workbook](#)

EM over PCFGs : an example

Discrete.
 $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

Obs/Symp.
 (y_1)

$y_1 = \text{"Mary saw a bird on a tree"}$

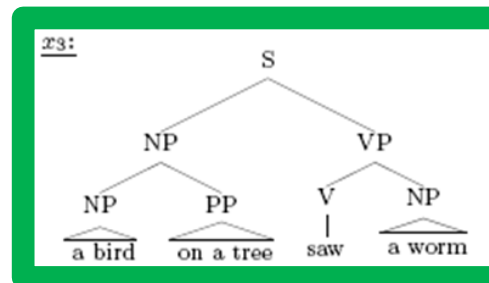


Inference

$$P(x_1 | y_1) = \frac{P(y_1 | x_1) P(x_1)}{P(y_1 | x_1) P(x_1) + P(y_1 | x_2) P(x_2)}$$

causing $P(x_2 | y_1)$

$y_2 = \text{"a bird on a tree saw a worm"}$

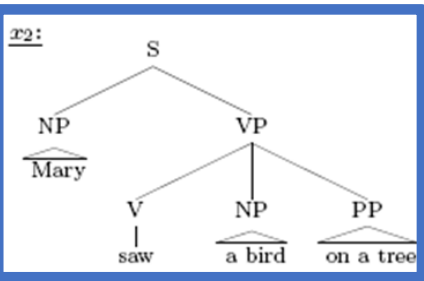
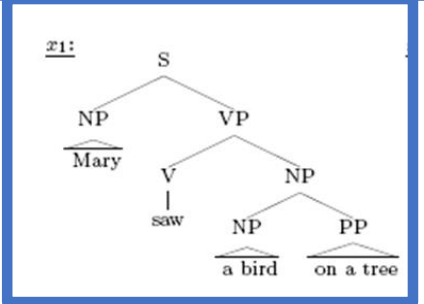


Refer Workbook

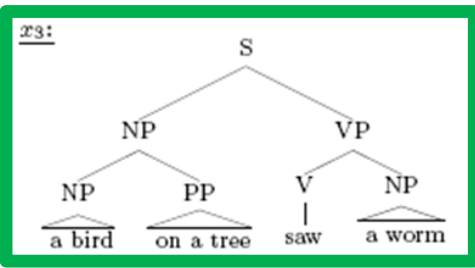
EM over PCFGs : an example (preparation for E step)

Preparatory :

$y_1 = \text{"Mary saw a bird on a tree"}$

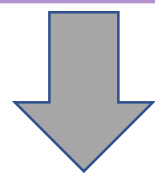


$y_2 = \text{"a bird on a tree saw a worm"}$



$$\begin{aligned} p_0(x_1) &= p(S \rightarrow NP VP) \cdot p\left(\frac{NP}{\text{Mary}}\right) \cdot p(VP \rightarrow V NP PP) \cdot p(V \rightarrow \text{saw}) \cdot p\left(\frac{NP}{\text{a bird}}\right) \cdot p\left(\frac{PP}{\text{on a tree}}\right) \\ &= 1.00 \cdot 0.25 \cdot 0.50 \cdot 1.00 \cdot 0.25 \cdot 0.25 \cdot 1.00 \\ &= 0.0078125 \end{aligned}$$

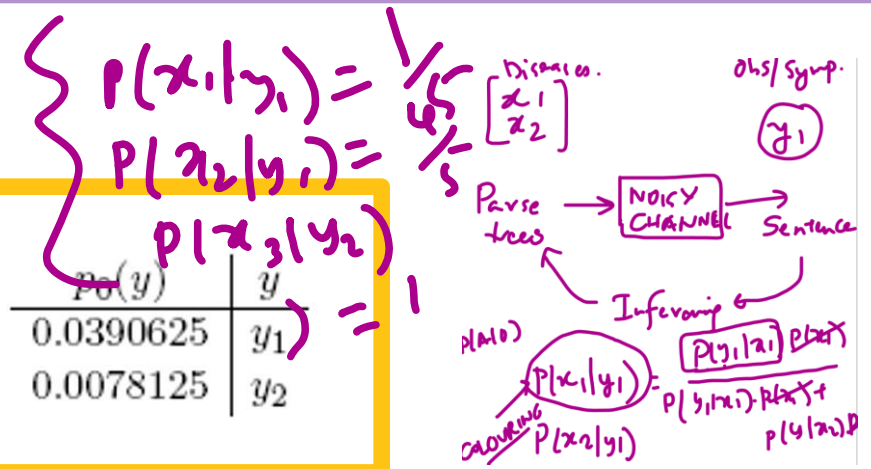
$p(y_1) = p_0(x_1) + p_0(x_2) = 0.0078125 + 0.0312500 = 0.0390625$



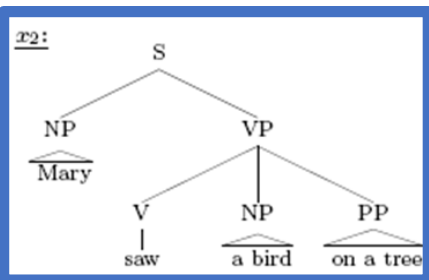
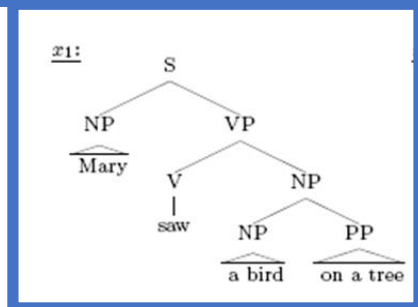
$p(y_1/x_1)$
 $p(y_1/x_2)$
 $p(y_2/x_3)$

$p_0(x)$	x
0.0078125	x_1 ✓ $\frac{1}{5}$
0.0312500	x_2 ✓
0.0078125	x_3 ✓

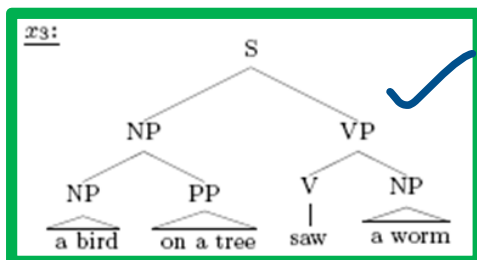
$p_0(y)$	y
0.0390625	y_1
0.0078125	y_2



$y_1 = \text{"Mary saw a bird on a tree"}$



$y_2 = \text{"a bird on a tree saw a worm"}$



$f(y)$	y	
5	y_1	$y_1 = \text{"Mary saw a bird on a tree"}$
10	y_2	$y_2 = \text{"a bird on a tree saw a worm"}$

$p_0(x)$	x
0.0078125	x_1
0.0312500	x_2
0.0078125	x_3

$p_0(y)$	y
0.0390625	y_1
0.0078125	y_2



$f_{T_q}(x)$	x
1	x_1
4	x_2
10	x_3

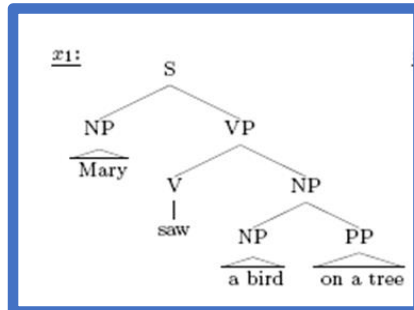
$$\begin{aligned}
 f_{T_q}(x_1) &= f(\text{yield}(x_1)) \cdot q(x_1 | \text{yield}(x_1)) \\
 &= f(y_1) \cdot q(x_1 | y_1) \\
 &= f(y_1) \cdot \frac{q(x_1)}{q(y_1)} \\
 &= 5 \cdot \frac{0.0078125}{0.0390625} \\
 &= 1
 \end{aligned}$$

Refer Workbook

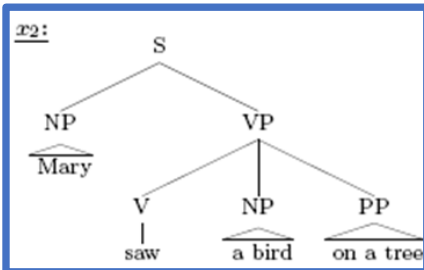
M step : estimate the PCFG parameters

$f_{T_q}(x)$	x
1	x_1
4	x_2
10	x_3

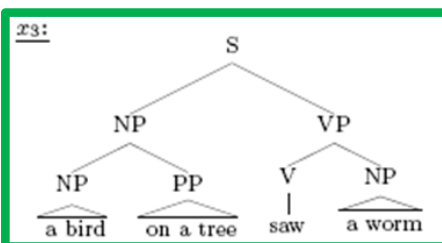
1 instance of



4 instances of



10 instances of



$S \longrightarrow NP VP \quad (1.000)$
 $VP \longrightarrow V NP \quad (0.733 \approx \frac{1+10}{15})$
 $VP \longrightarrow V NP PP \quad (0.267 \approx \frac{4}{15})$
 $NP \longrightarrow NP PP \quad (0.268 \approx \frac{1+10}{41})$
 $NP \longrightarrow Mary \quad (0.122 \approx \frac{1+4}{41})$
 $NP \longrightarrow a \text{ bird} \quad (0.366 \approx \frac{1+4+10}{41})$
 $NP \longrightarrow a \text{ worm} \quad (0.244 \approx \frac{10}{41})$
 $PP \longrightarrow on \text{ a tree} \quad (1.000)$
 $V \longrightarrow saw \quad (1.000)$

[Refer Workbook](#)

M step: comparison of results of M_1 and M_0 steps

After M_1 step

After M_0 step

$S \rightarrow NP VP$	(1.00)
$VP \rightarrow V NP$	(0.50)
$VP \rightarrow V NP PP$	(0.50)
$NP \rightarrow NP PP$	(0.25)
$NP \rightarrow Mary$	(0.25)
$NP \rightarrow a bird$	(0.25)
$NP \rightarrow a worm$	(0.25)
$PP \rightarrow on a tree$	(1.00)
$V \rightarrow saw$	(1.00)

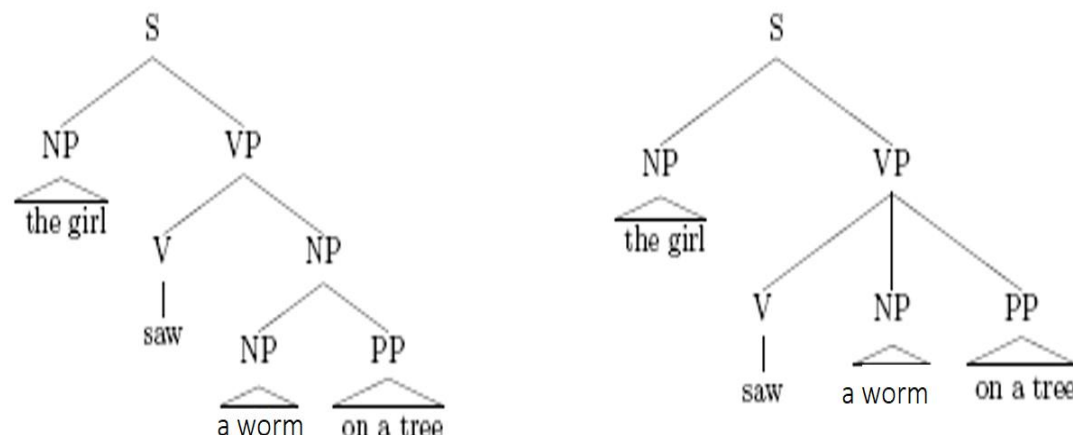
$S \rightarrow NP VP$	(1.000)
$VP \rightarrow V NP$	$(0.733 \approx \frac{1+10}{15})$
$VP \rightarrow V NP PP$	$(0.267 \approx \frac{4}{15})$
$NP \rightarrow NP PP$	$(0.268 \approx \frac{1+10}{41})$
$NP \rightarrow Mary$	$(0.122 \approx \frac{1+4}{41})$
$NP \rightarrow a bird$	$(0.366 \approx \frac{1+4+10}{41})$
$NP \rightarrow a worm$	$(0.244 \approx \frac{10}{41})$
$PP \rightarrow on a tree$	(1.000)
$V \rightarrow saw$	(1.000)

Finally ...

CFG rule	p_0	p_1	p_2	p_3	...	p_{18}
S \longrightarrow NP VP	1.000	1.000	1.000	1.000		1.000
VP \longrightarrow V NP	0.500	0.733	0.807	0.850		0.967
VP \longrightarrow V NP PP	0.500	0.267	0.193	0.150		0.033
NP \longrightarrow NP PP	0.250	0.268	0.287	0.298		0.326
NP \longrightarrow Mary	0.250	0.122	0.118	0.117		0.112
NP \longrightarrow a bird	0.250	0.366	0.357	0.351		0.337
NP \longrightarrow a worm	0.250	0.244	0.238	0.234		0.225
PP \longrightarrow on a tree	1.000	1.000	1.000	1.000		1.000
V \longrightarrow saw	1.000	1.000	1.000	1.000		1.000

[Refer Workbook](#)

Using the PCFG for disambiguation



$$p(\text{VP} \rightarrow \text{V NP}) \cdot p(\text{NP} \rightarrow \text{NP PP}) > p(\text{VP} \rightarrow \text{V NP PP})$$

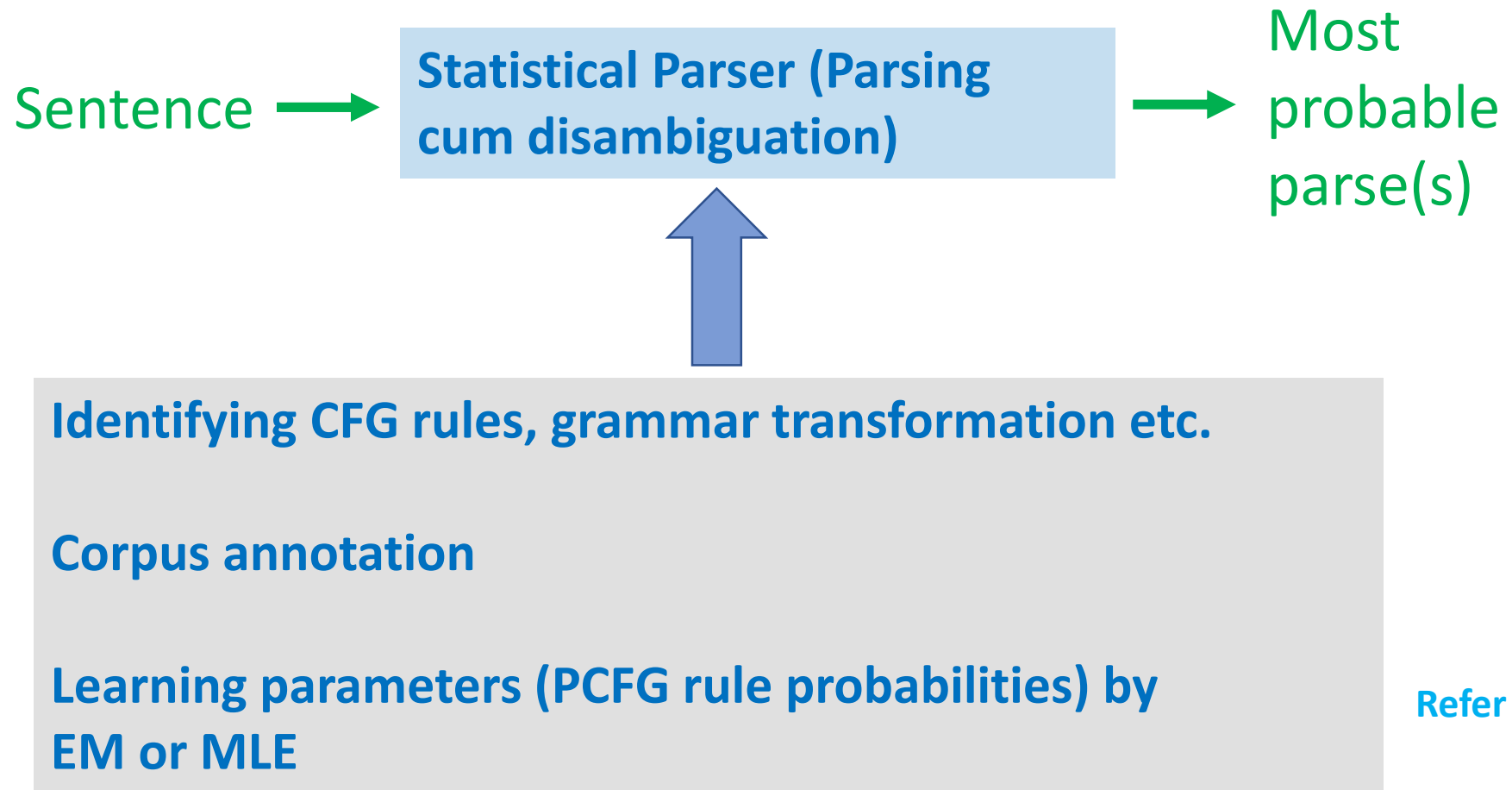
CFG rule	p_0	p_1	p_2	p_3	...	p_{18}
$S \rightarrow \text{NP VP}$	1.000	1.000	1.000	1.000		1.000
$\text{VP} \rightarrow \text{V NP}$	0.500	0.733	0.807	0.850		0.967
$\text{VP} \rightarrow \text{V NP PP}$	0.500	0.267	0.193	0.150		0.033
$\text{NP} \rightarrow \text{NP PP}$	0.250	0.268	0.287	0.298		0.326
$\text{NP} \rightarrow \text{Mary}$	0.250	0.122	0.118	0.117		0.112
$\text{NP} \rightarrow \text{a bird}$	0.250	0.366	0.357	0.351		0.337
$\text{NP} \rightarrow \text{a worm}$	0.250	0.244	0.238	0.234		0.225
$\text{PP} \rightarrow \text{on a tree}$	1.000	1.000	1.000	1.000		1.000
$\text{V} \rightarrow \text{saw}$	1.000	1.000	1.000	1.000		1.000

The parse on the left is preferred if:

$$p(\text{VP} \rightarrow \text{V NP}) \cdot p(\text{NP} \rightarrow \text{NP PP}) > p(\text{VP} \rightarrow \text{V NP PP})$$

p	$p(\text{VP} \rightarrow \text{V NP}) \cdot p(\text{NP} \rightarrow \text{NP PP})$	$p(\text{VP} \rightarrow \text{V NP PP})$
p_0	$0.500 \cdot 0.250 = 0.125$	0.500
p_1	$0.733 \cdot 0.268 = 0.196$	0.267
p_2	$0.807 \cdot 0.287 = \mathbf{0.232}$	0.193
p_3	$0.850 \cdot 0.298 = \mathbf{0.253}$	0.150
\vdots		
p_{18}	$0.967 \cdot 0.326 = \mathbf{0.315}$	0.033

PCFGs : The Big Picture



Refer Workbook

Perpetual Motion Machine?

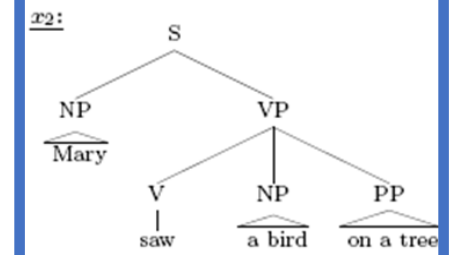
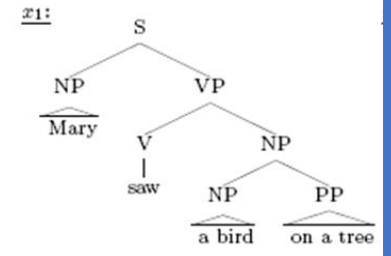
- Is this magic? Where does the extra knowledge come from?



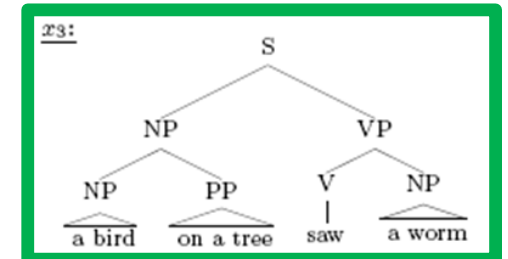
Ack:Wikipedia

Perpetual motion wheels from a drawing by [Leonardo da Vinci](#)

$y_1 = \text{"Mary saw a bird on a tree"}$



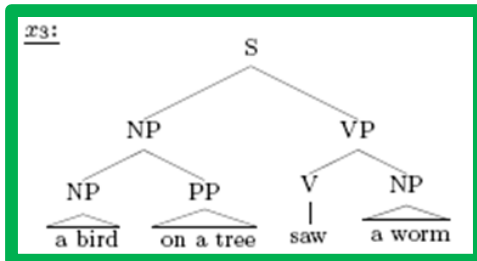
$y_2 = \text{"a bird on a tree saw a worm"}$



Perpetual Motion Machine?

- Is this magic? Where does the extra knowledge come from?

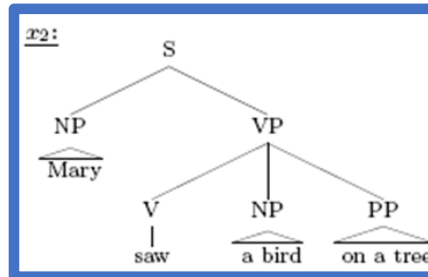
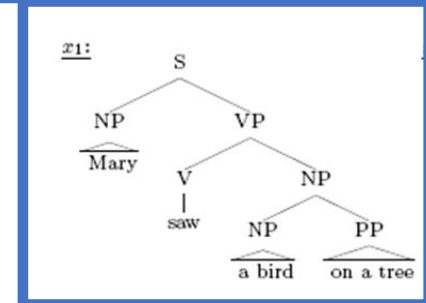
$y_2 = \text{"a bird on a tree saw a worm"}$



Ack:Wikipedia

helps
disambiguate

$y_1 = \text{"Mary saw a bird on a tree"}$



The Unifying Picture

	K-means	EM	Biased Coins	PCFG
Source				
Observations				
Generative storyline (preparation for E step)				
Parameters Estimated				
Why it works				

Reference

A Tutorial on the Expectation-Maximization Algorithm
Including Maximum-Likelihood Estimation and EM Training of
Probabilistic Context-Free Grammars

By Detlef Prescher