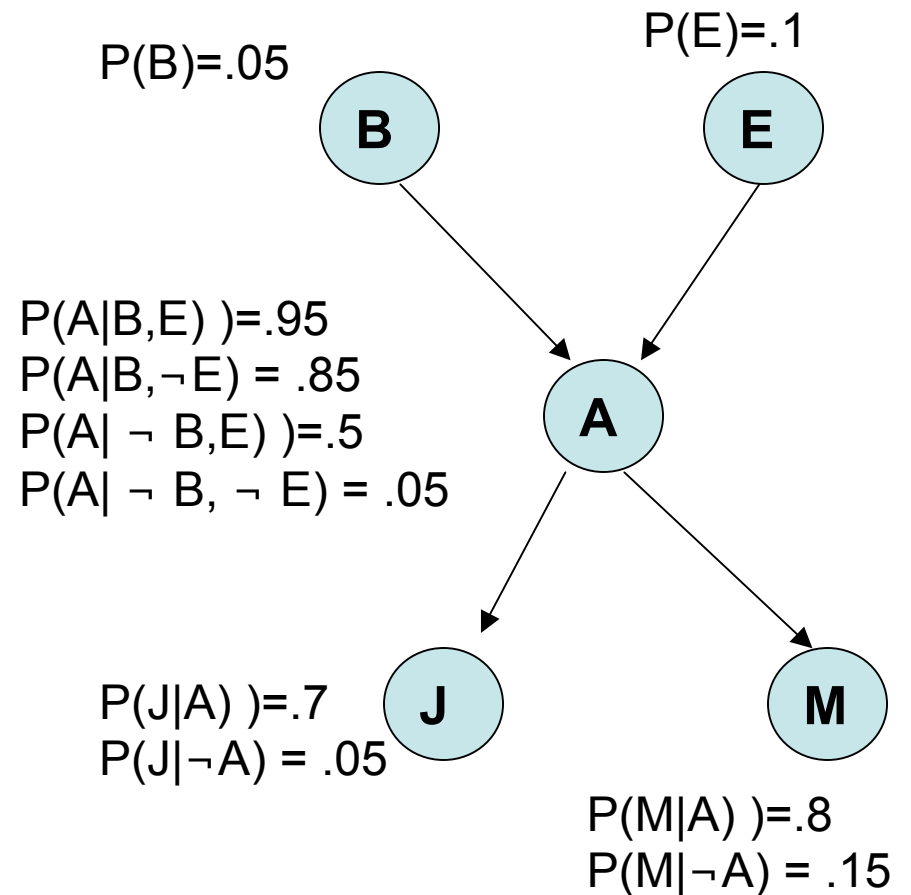


15-780: Graduate Artificial Intelligence

Density estimation

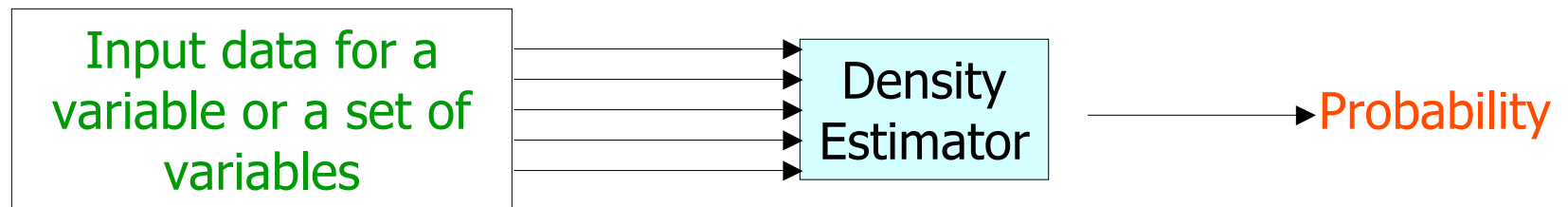
Conditional Probability Tables (CPT)

But where do we get them?



Density Estimation

- A Density Estimator learns a mapping from a set of attributes to a Probability



Density estimation

- Estimate the distribution (or conditional distribution) of a random variable
- Types of variables:
 - Binary
coin flip, alarm
 - Discrete
dice, car model year
 - Continuous
height, weight, temp.,

Not just for Bayesian networks ...

- Density estimators can do many good things...
 - Can sort the records by probability, and thus spot weird records (anomaly detection)
 - Can do inference: $P(E1|E2)$
 - Medical diagnosis / Robot sensors
 - Ingredient for Bayes networks

Density estimation

- Binary and discrete variables:

Easy: Just count!

- Continuous variables:

Harder (but just a bit): Fit a model

Learning a density estimator

$$\hat{P}(x[i] = u) = \frac{\text{\#records in which } x[i] = u}{\text{total number of records}}$$

A trivial learning algorithm!

Course evaluation

$$P(\text{summer}) = \# \text{Summer} / \# \text{ records} \\ = 23/151 = 0.15$$

$$P(\text{Evaluation} = 1) = \# \text{Evaluation}=1 \\ / \# \text{ records} \\ = 49/151 = 0.32$$

$$P(\text{Evaluation} = 1 \mid \text{summer}) = \\ P(\text{Evaluation} = 1 \ \& \ \text{summer}) / \\ P(\text{summer}) = 2/23 = 0.09$$

But why do we count?

Summer?	Size	Evaluation
1	19	3
1	17	3
0	49	2
0	33	1
0	55	3
1	20	1



Computing the joint likelihood of the data

$P(\text{summer}) = \# \text{Summer} / \# \text{ records}$
 $= 23/151 = 0.15$

Summer?	Size	Evaluation

$$\hat{P}(\text{dataset}|M) = \hat{P}(\mathbf{x}_1 \wedge \mathbf{x}_2 \dots \wedge \mathbf{x}_R|M) = \prod_{k=1}^R \hat{P}(\mathbf{x}_k|M)$$

The next slide presents one of the most important ideas in probabilistic inference. It has a huge number of applications in many different and diverse problems

Maximum Likelihood Principle

- We can fit models by maximizing the probability of generating the observed samples:

$$L(x_1, \dots, x_n | \Theta) = p(x_1 | \Theta) \dots p(x_n | \Theta)$$

- The samples (rows in the table) are assumed to be independent)
- For a binary random variable A with $P(A=1)=q$
 $\operatorname{argmax}_q = \#1/\#\text{samples}$
- Why?

Maximum Likelihood Principle

- For a binary random variable A with $P(A=1)=q$
 $\operatorname{argmax}_q = \#1/\#\text{samples}$
- Why?

Data likelihood: $P(D | M) = q^{n_1} (1 - q)^{n_2}$

We would like to find: $\operatorname{argmax}_q q^{n_1} (1 - q)^{n_2}$

Maximum Likelihood Principle

Data likelihood: $P(D | M) = q^{n_1} (1 - q)^{n_2}$

We would like to find: $\arg \max_q q^{n_1} (1 - q)^{n_2}$

$$\frac{\partial}{\partial q} q^{n_1} (1 - q)^{n_2} = n_1 q^{n_1 - 1} (1 - q)^{n_2} - q^{n_1} n_2 (1 - q)^{n_2 - 1}$$

$$\frac{\partial}{\partial q} = 0 \Rightarrow$$

$$n_1 q^{n_1 - 1} (1 - q)^{n_2} - q^{n_1} n_2 (1 - q)^{n_2 - 1} = 0 \Rightarrow$$

$$q^{n_1 - 1} (1 - q)^{n_2 - 1} (n_1 (1 - q) - q n_2) = 0 \Rightarrow$$

$$n_1 (1 - q) - q n_2 = 0 \Rightarrow$$

$$n_1 = n_1 q + n_2 q \Rightarrow$$

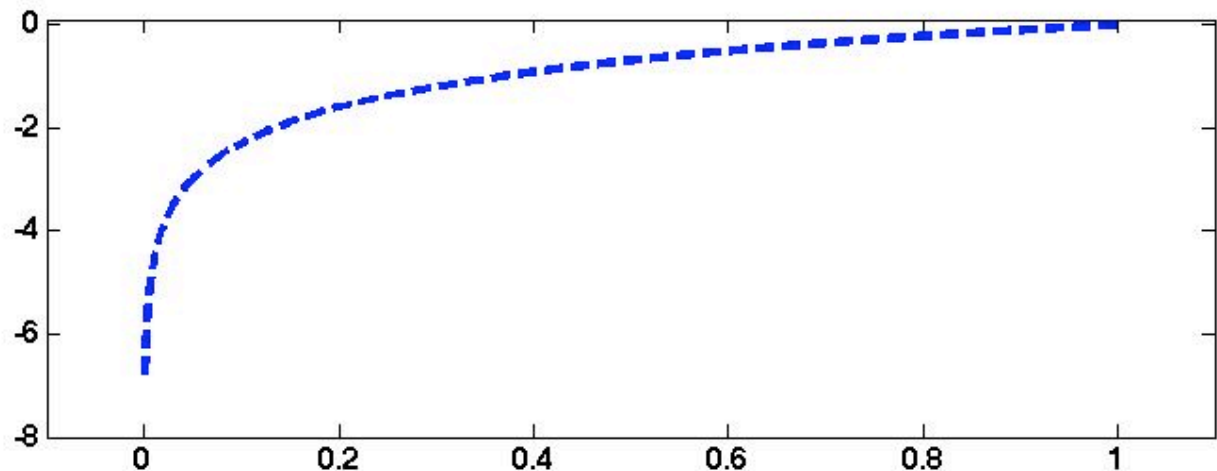
$$q = \frac{n_1}{n_1 + n_2}$$

Log Probabilities

When working with products, probabilities of entire datasets often get too small. A possible solution is to use the log of probabilities, often termed 'log likelihood'

$$\log \hat{P}(\text{dataset}|M) = \log \prod_{k=1}^R \hat{P}(\mathbf{x}_k|M) = \sum_{k=1}^R \log \hat{P}(\mathbf{x}_k|M)$$

Log values
between 0 and 1



Density estimation

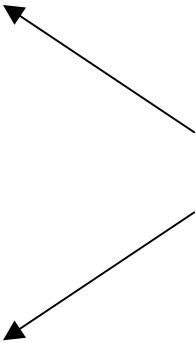
- Binary and discrete variables:

Easy: Just count!

- Continuous variables:

Harder (but just a bit): Fit a model

But what if we only have very few samples?



The danger of joint density estimation

$P(\text{summer} \ \& \ \text{size} > 20 \ \& \ \text{evaluation} = 3)$
 $= 0$

- No such example in our dataset

Now lets assume we are given a new (often called 'test') dataset. If this dataset contains the line

Summer	Size	Evaluation
1	30	3

Then the probability we would assign to the *entire* dataset is 0

Summer?	Size	Evaluation
1	19	3
1	17	3
0	49	2
0	33	1
0	55	3
1	20	1

Naïve Density Estimation

The problem with the Joint Estimator is that it just mirrors the training data.

We need something which generalizes more usefully.

The **naïve model** generalizes strongly:

Assume that each attribute is distributed independently of any of the other attributes.

Joint estimation, revisited

Assuming independence we can compute each probability independently

$$P(\text{Summer}) = 0.15$$

$$P(\text{Evaluation} = 1) = 0.32$$

$$P(\text{Size} > 20) = 0.63$$

How do we do on the joint?

$$P(\text{Summer} \& \text{Evaluation} = 1) = 0.09$$

$$P(\text{Summer})P(\text{Evaluation} = 1) = 0.05$$

$$P(\text{size} > 20 \& \text{Evaluation} = 1) = 0.23$$

$$P(\text{size} > 20)P(\text{Evaluation} = 1) = 0.20$$

Summer?	Size	Evaluation
1	19	3
1	17	3
0	49	2
		1
		3
1	20	1

Not bad !

Joint estimation, revisited

Assuming independence we can compute each probability independently

$$P(\text{Summer}) = 0.15$$

$$P(\text{Evaluation} = 1) = 0.32$$

$$P(\text{Size} > 20) = 0.63$$

How do we do on the joint?

$$P(\text{Summer} \& \text{Size} > 20) = 0.026$$

$$P(\text{Summer})P(\text{Size} > 20) = 0.094$$

Summer?	Size	Evaluation
1	19	3
1	17	3
0	49	2
0	33	1
0	55	3
1	20	1

We must be careful when using the Naïve density estimator

Contrast

Joint DE	Naïve DE
Can model anything	Can model only very boring distributions
No problem to model “C is a noisy copy of A”	Outside Naïve’s scope
Given 100 records and more than 6 Boolean attributes will screw up badly	Given 100 records and 10,000 multivalued attributes will be fine

Dealing with small datasets

- We just discussed one possibility: Naïve estimation
- There is another way to deal with small number of measurements that is often used in practice.
- Assume we want to compute the probability of heads in a coin flip
 - What if we can only observe 3 flips?
 - 25% of the times a maximum likelihood estimator will assign probability of 1 to either the heads or tails



Pseudo counts

- What if we can only observe 3 flips?
- 25% of the times a maximum likelihood estimator will assign probability of 1 to either the heads or tails
- In these cases we can use prior belief about the 'fairness' of most coins to influence the resulting model.
- We assume that we have observed 10 flips with 5 tails and 5 heads
- Thus $p(\text{heads}) = (\#\text{heads}+5)/(\#\text{flips}+10)$
- Advantages: 1. Never assign a probability of 0 to an event
2. As more data accumulates we can get very close to the real distribution (the impact of the pseudo counts will diminish rapidly)

Pseudo counts

- What if we can only observe 3 flips?
- 25% of the times a maximum likelihood estimator will assign probability of 1 to either the head or tail

- In these cases, the maximum likelihood estimator is not a good estimator of the true probability of heads, because it is biased and does not reflect the 'fairness' of the coin.

model.

- We assume the coin is fair and flip it 5 times.

Some distributions (for example, the Beta distribution) can incorporate pseudo counts as part of the model

5 tails

- Thus, the maximum likelihood estimator is not a good estimator of the true probability of heads.

- Advantages of using pseudo counts:

1. It is a better estimator of the true probability of heads.

- 2. As more data is collected, the maximum likelihood estimator converges to the true probability of heads.

real

distribution (the impact of the pseudo counts will diminish rapidly)

Density estimation

- Binary and discrete variables:

Easy: Just count!

✓

- Continuous variables:

Harder (but just a bit): Fit a model

Conditional Probability Tables (CPT)

**What do we do with
continuous variables?**

S1 – sensor 1

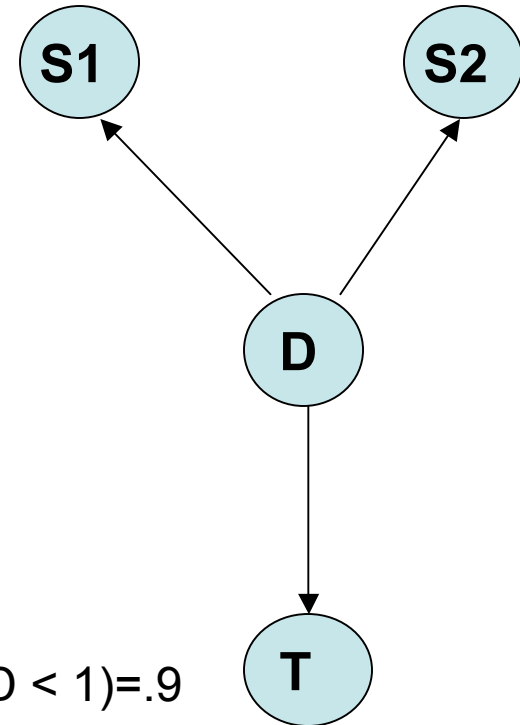
S2 – sensor 2

D – distance to wall

T – too close

$$P(S1 | D) = ?$$

$$P(S2 | D) = ?$$



$$P(T | D < 1) = .9$$

Conditional Probability Tables (CPT)

**What do we do with
continuous variables?**

S1 – sensor 1

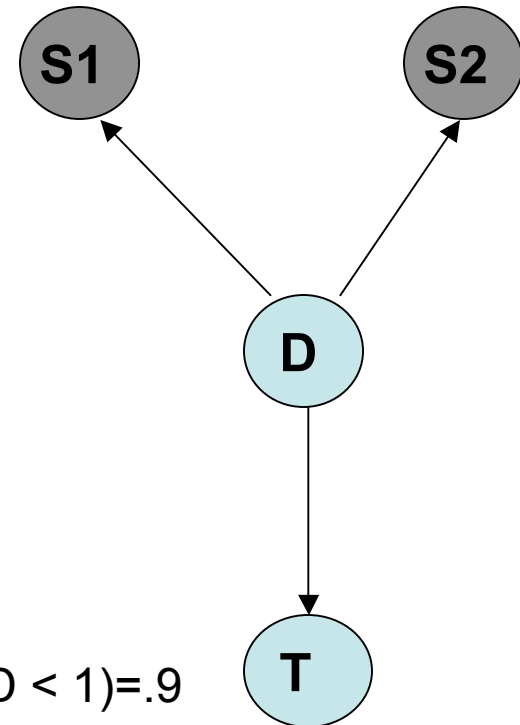
S2 – sensor 2

D – distance to wall

T – too close

$$P(S1 | D) = ?$$

$$P(S2 | D) = ?$$



$$P(T | D < 1) = .9$$

Elementary Concepts

- **Population:** the ideal group whose properties we are interested in and from which the samples are drawn
e.g., graduate students at CMU
- **Random sample:** a set of elements drawn at random from the population
e.g., students in grad AI

Elementary Concepts

- **Statistic:** a number computed from the data
e.g., Average time of sleep

Sample Statistics

- **Sample mean:** $\overline{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$

where n is the number of samples.

- **Sample variance:**

$$\overline{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (x_i - \overline{\mu})^2$$

- **Sample covariance:**

$$\overline{\text{COV}(x_1, x_2)} = \frac{1}{n} \sum_{i=1}^n (x_{1,i} - \overline{\mu_1}) (x_{2,i} - \overline{\mu_2})$$

How much do grad students sleep?

- Lets try to estimate the distribution of the time graduate students spend sleeping (outside class).

Possible statistics

- **X**

Sleep time

- **Mean of X:**

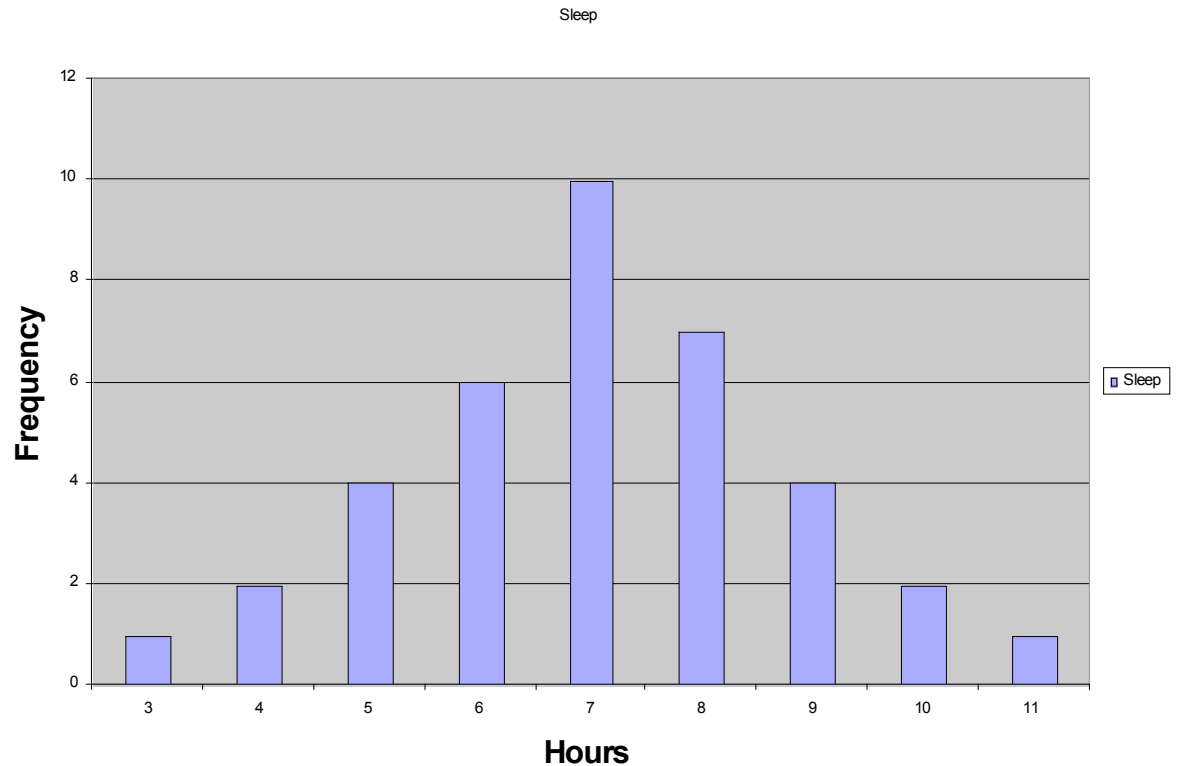
$$E\{X\}$$

7.03

- **Variance of X:**

$$\text{Var}\{X\} = E\{(X - E\{X\})^2\}$$

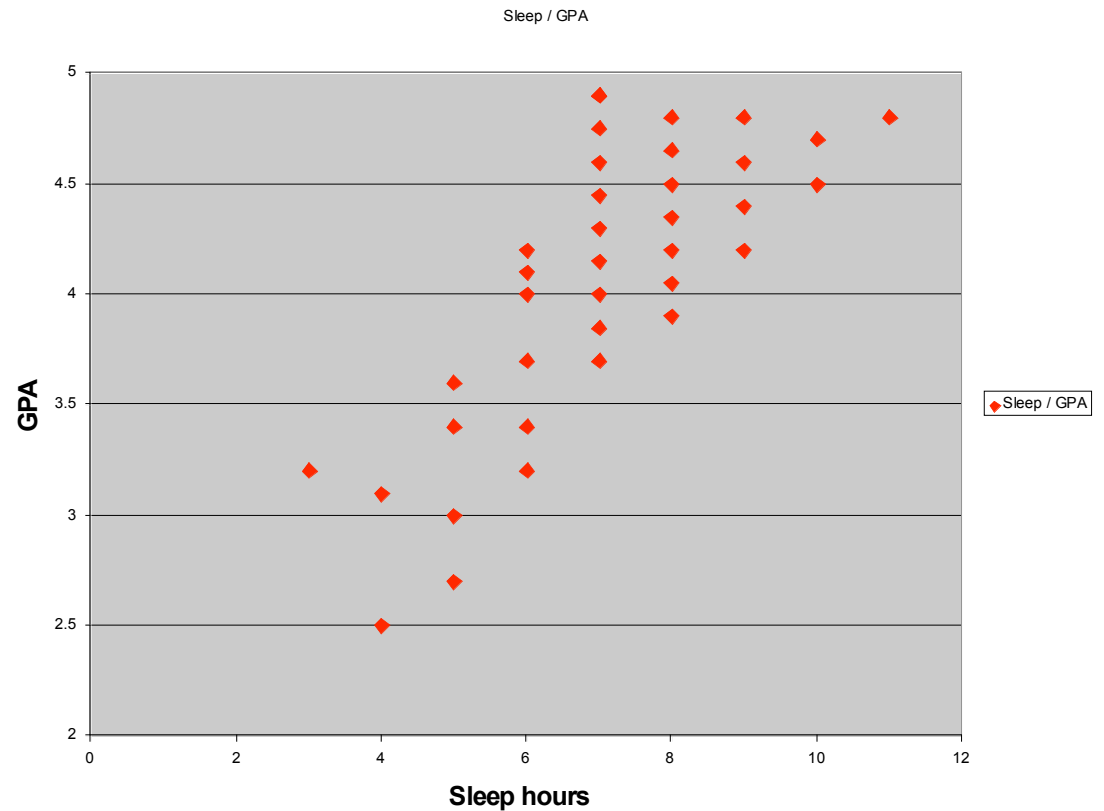
3.05



Covariance: Sleep vs. GPA

•Co-Variance of X1,
X2:

$$\begin{aligned} \text{Covariance}\{X1, X2\} &= \\ E\{(X1 - E\{X1\})(X2 - E\{X2\})\} &= \\ &= 0.88 \end{aligned}$$



Statistical Models

- Statistical models attempt to characterize properties of the population of interest
- For example, we might believe that repeated measurements follow a normal (Gaussian) distribution with some mean μ and variance σ^2 , $x \sim N(\mu, \sigma^2)$

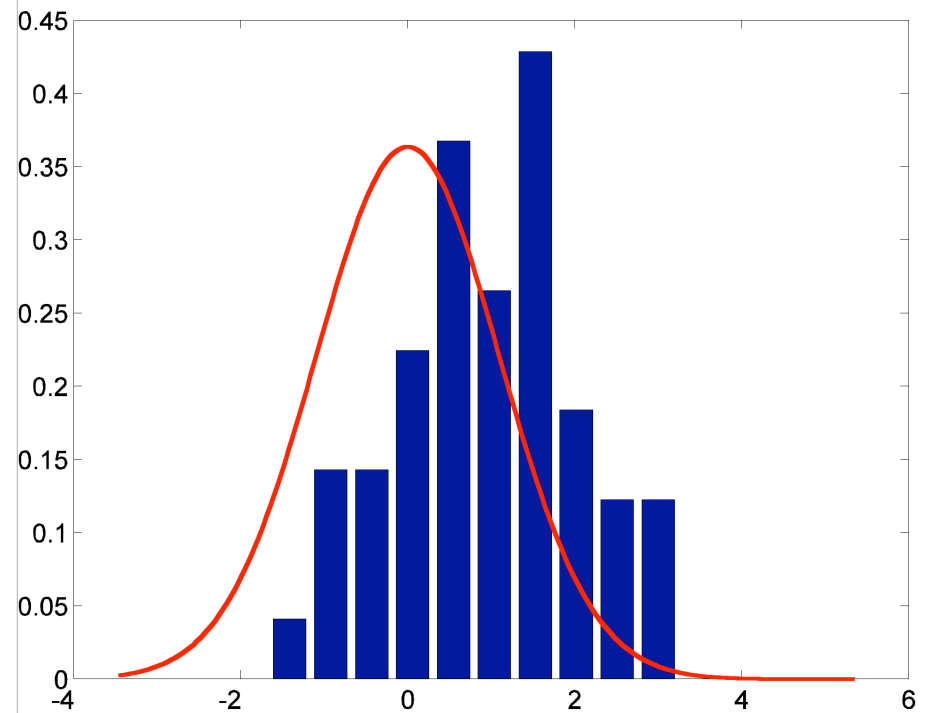
where

$$p(x | \Theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

and $\Theta=(\mu, \sigma^2)$ defines the parameters (mean and variance) of the model.

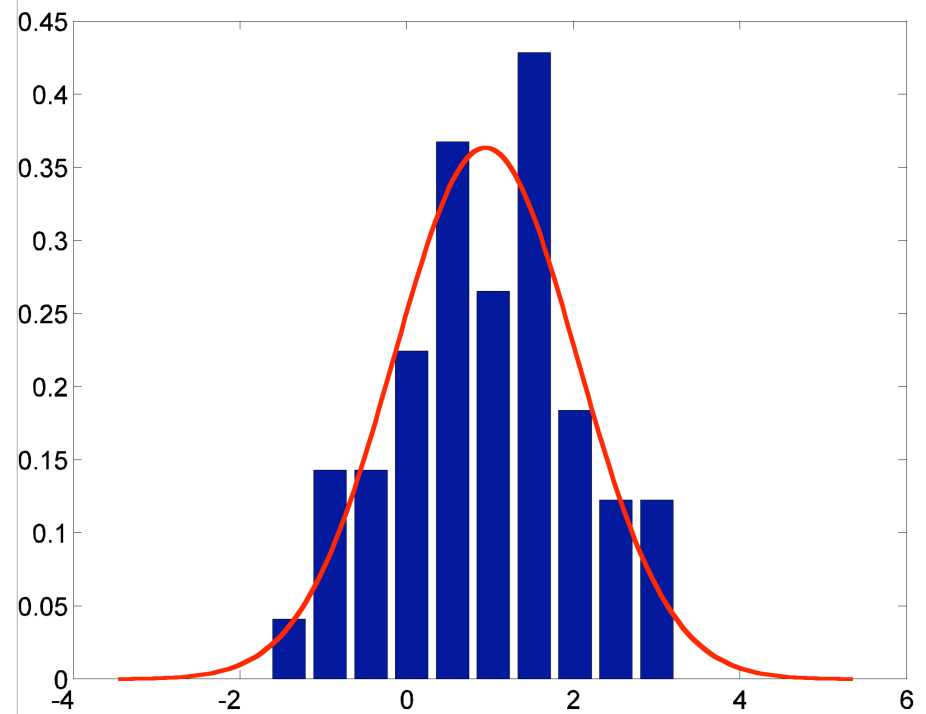
The Parameters of Our Model

- A statistical model is a **collection** of distributions; the **parameters** specify individual distributions $x \sim N(\mu, \sigma^2)$
- We need to adjust the parameters so that the resulting distribution **fits** the data well



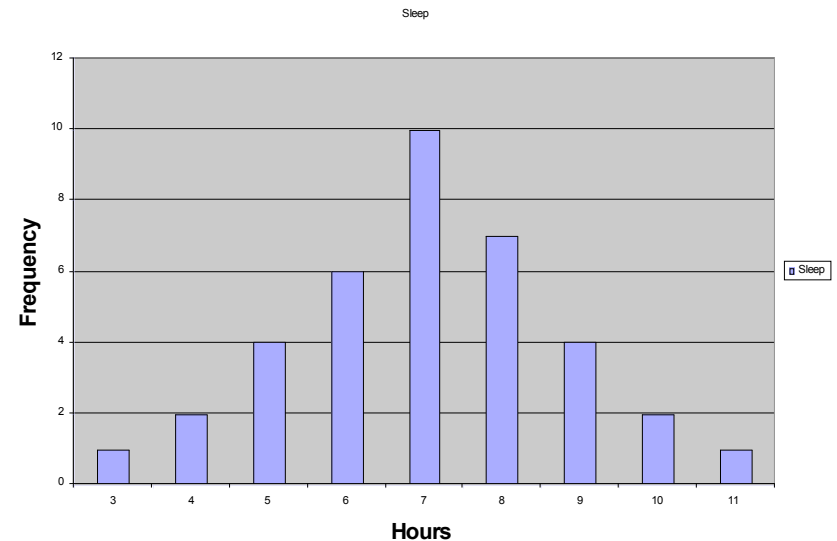
The Parameters of Our Model

- A statistical model is a **collection** of distributions; the **parameters** specify individual distributions $x \sim N(\mu, \sigma^2)$
- We need to adjust the parameters so that the resulting distribution **fits** the data well



Computing the parameters of our model

- Lets assume a Gaussian distribution for our sleep data
- How do we compute the parameters of the model?



Maximum Likelihood Principle

- We can fit statistical models by maximizing the probability of generating the observed samples:

$$L(x_1, \dots, x_n | \Theta) = p(x_1 | \Theta) \dots p(x_n | \Theta)$$

(the samples are assumed to be independent)

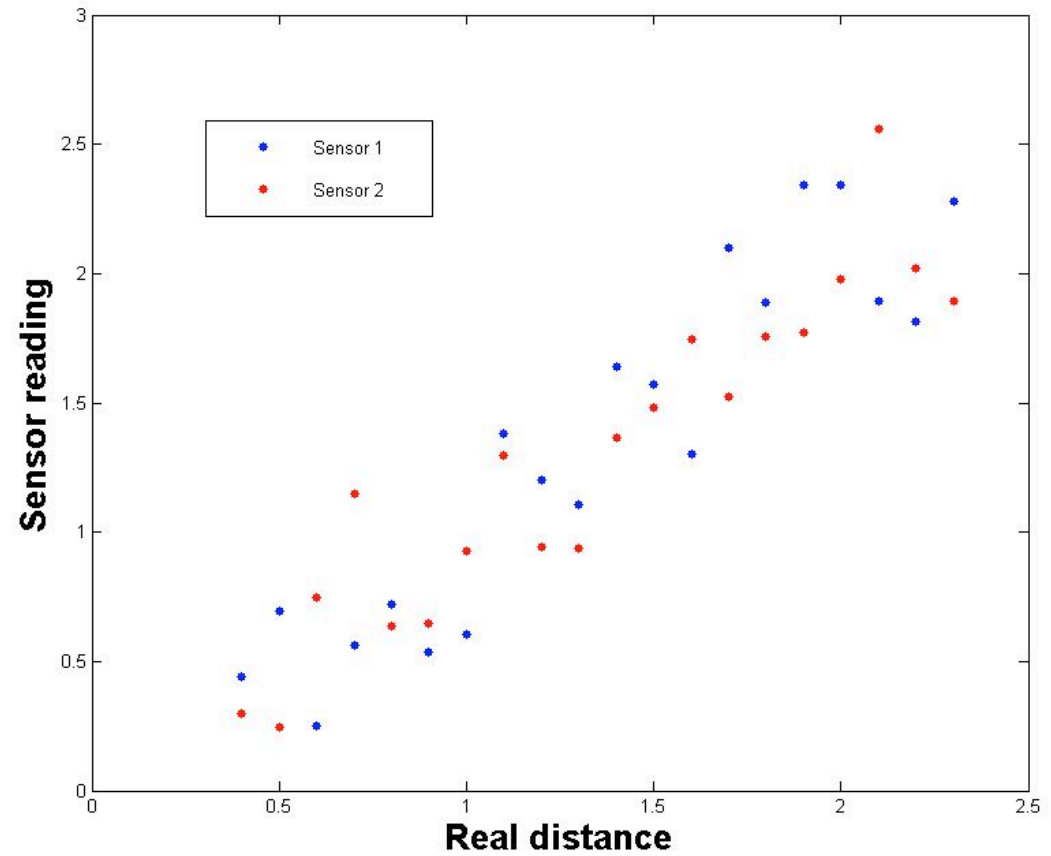
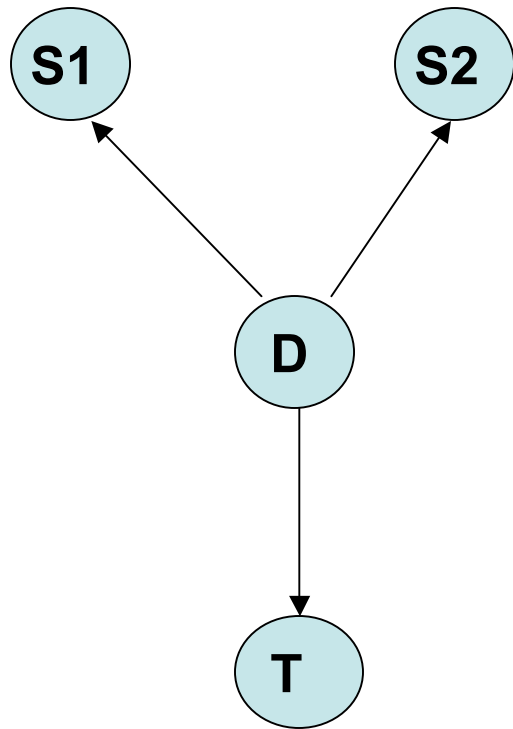
- In the Gaussian case we simply set the mean and the variance to the sample mean and the sample variance:

$$\bar{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{\mu})^2$$

Why?

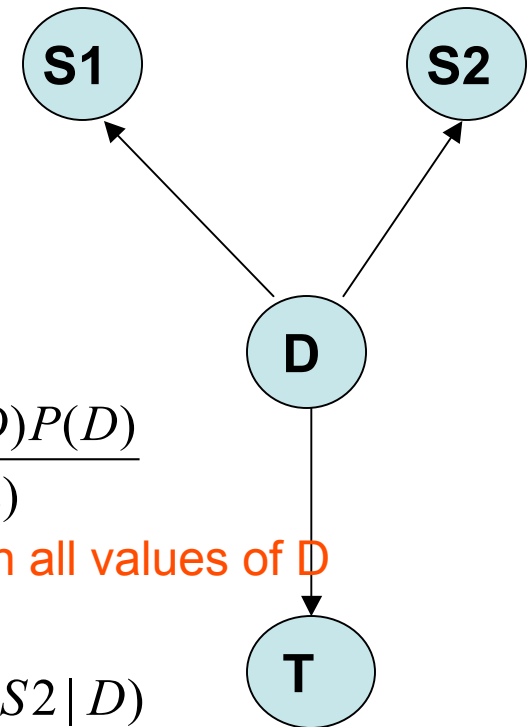
I will leave these derivation to you ...

Sensor data



What value would we infer for D given S1, S2?

- We will write the general terms and then use the network model to simplify it.
- The important issue is how to work with Gaussians



$$P(D | S1, S2) = \frac{P(S1 | D, S2)P(D | S2)}{P(S1 | S2)} \stackrel{\text{Bayes rule}}{=} \frac{P(S1 | D, S2)P(S2 | D)P(D)}{P(S1 | S2)P(S2)}$$

$$\stackrel{\text{Using network structure}}{\arg \max_D} \frac{P(S1 | D)P(S2 | D)P(D)}{P(S1 | S2)P(S2)} \stackrel{\text{Assuming equal prior on all values of D}}{=} \arg \max_D P(S1 | D)P(S2 | D)$$

$$P(S1 | D)P(S2 | D) = \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(D-S1)^2}{2\sigma_1^2}} \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{(D-S2)^2}{2\sigma_2^2}}$$

Model for sensor data

$$\log\left(\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(D-S1)^2}{2\sigma_1^2}} \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{(D-S2)^2}{2\sigma_2^2}}\right) = \log\left(\frac{1}{\sqrt{2\pi\sigma_1^2}} \frac{1}{\sqrt{2\pi\sigma_2^2}}\right) - \frac{(D-S1)^2}{2\sigma_1^2} - \frac{(D-S2)^2}{2\sigma_2^2}$$

$$\frac{\partial}{\partial D} \log\left(\frac{1}{\sqrt{2\pi\sigma_1^2}} \frac{1}{\sqrt{2\pi\sigma_2^2}}\right) - \frac{(D-S1)^2}{2\sigma_1^2} - \frac{(D-S2)^2}{2\sigma_2^2} = -2 \frac{(D-S1)}{2\sigma_1^2} - 2 \frac{(D-S2)}{2\sigma_2^2}$$

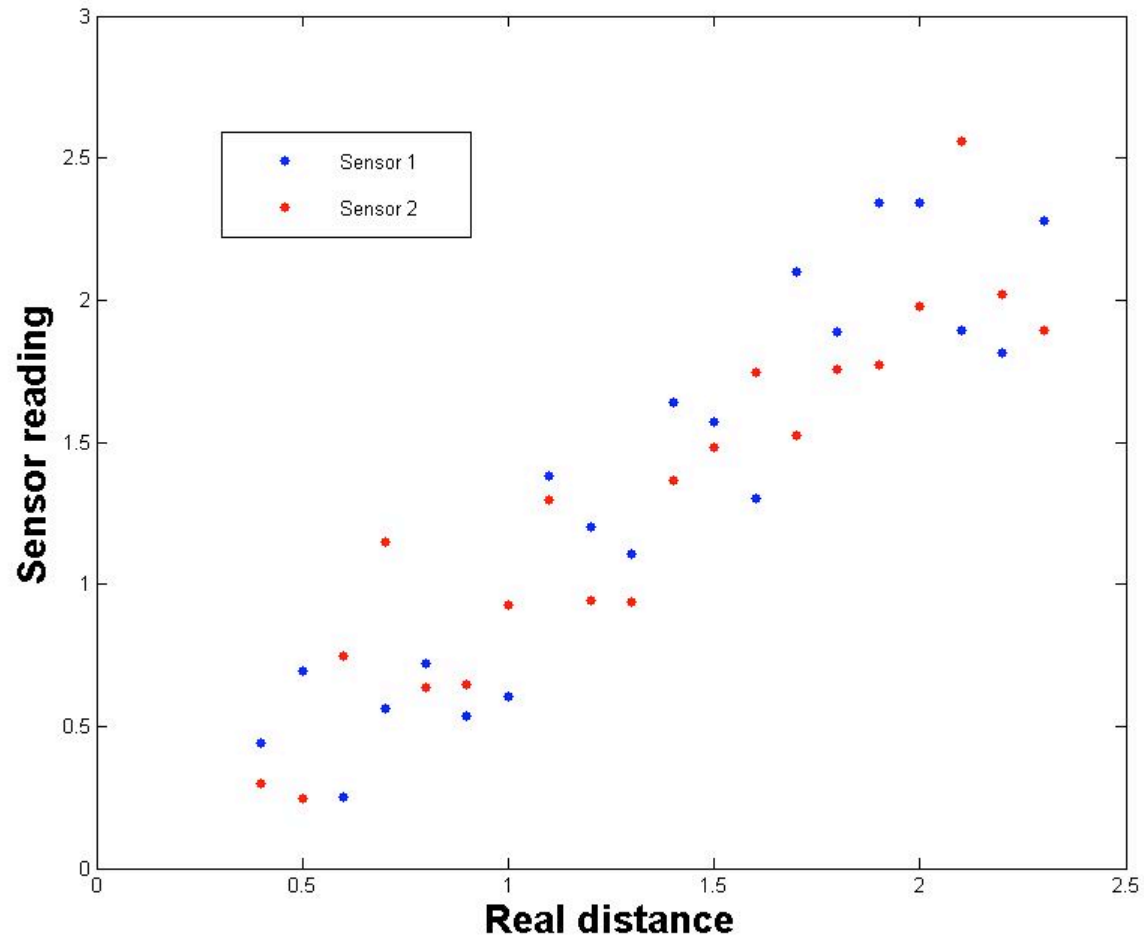
$$\Rightarrow -2 \frac{(D-S1)}{2\sigma_1^2} - 2 \frac{(D-S2)}{2\sigma_2^2} = 0 \Rightarrow$$

$$D = \frac{S1\sigma_2^2 + S2\sigma_1^2}{\sigma_1^2 + \sigma_2^2} \Rightarrow$$

$$D = \frac{S1 + S2}{2}$$

Only if $\sigma_1 = \sigma_2$

Sensor data



$$D = \frac{S1 + S2}{2}$$

Lets go back to Naïve vs.full model

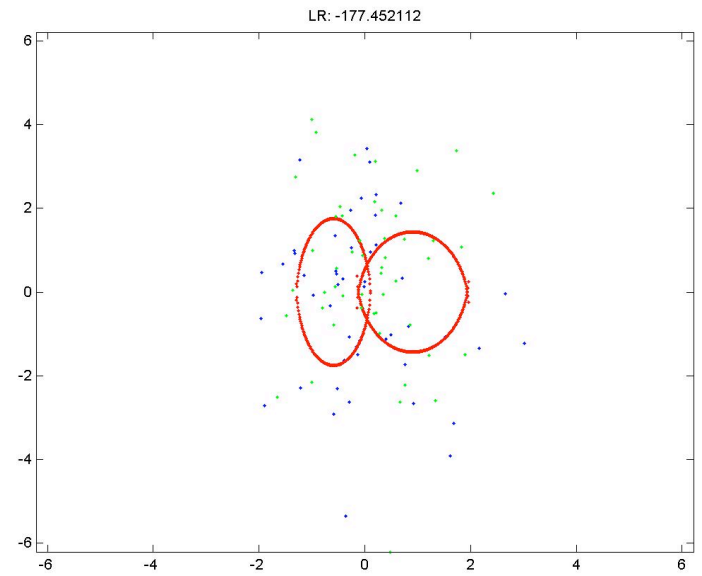
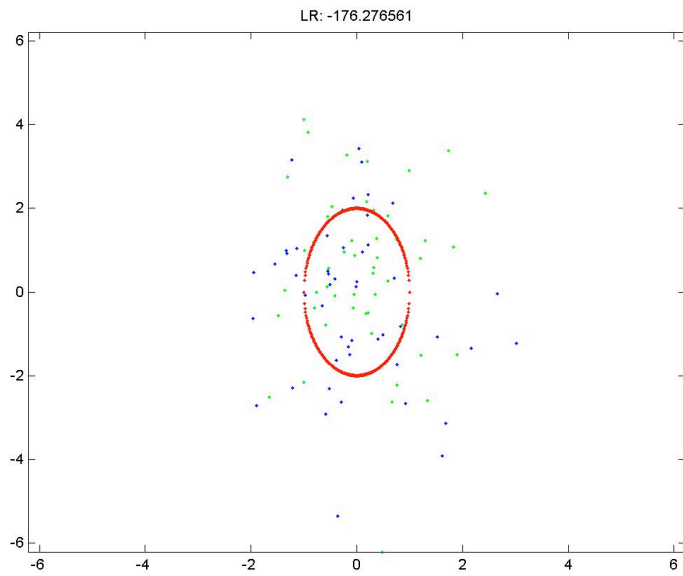
What should I use?

This can be determined based on:

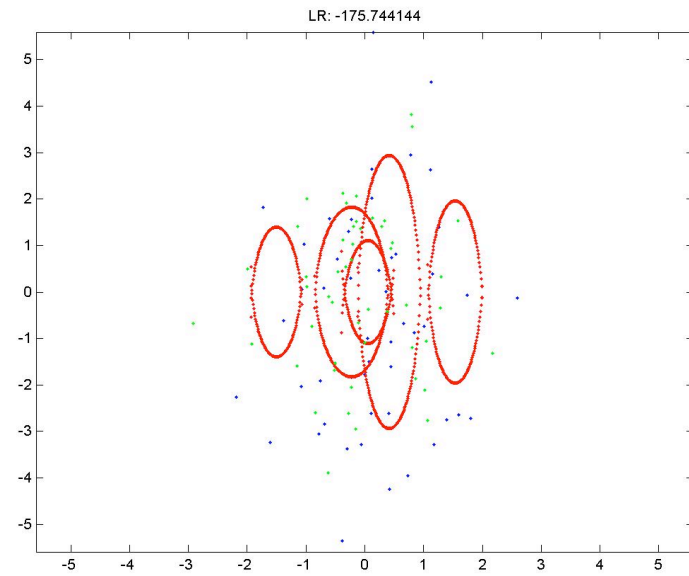
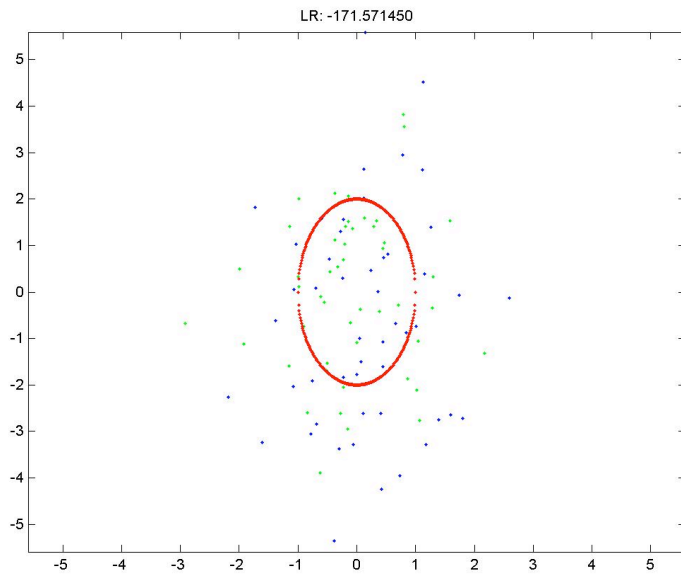
- Training data size
- Cross validation
- Likelihood ratio test

Cross validation is one of
the most useful tricks in
model fitting

Cross validation



Cross validation



Multi-Variate Gaussian

- A multivariate Gaussian model: $\mathbf{x} \sim N(\mu, \Sigma)$ where

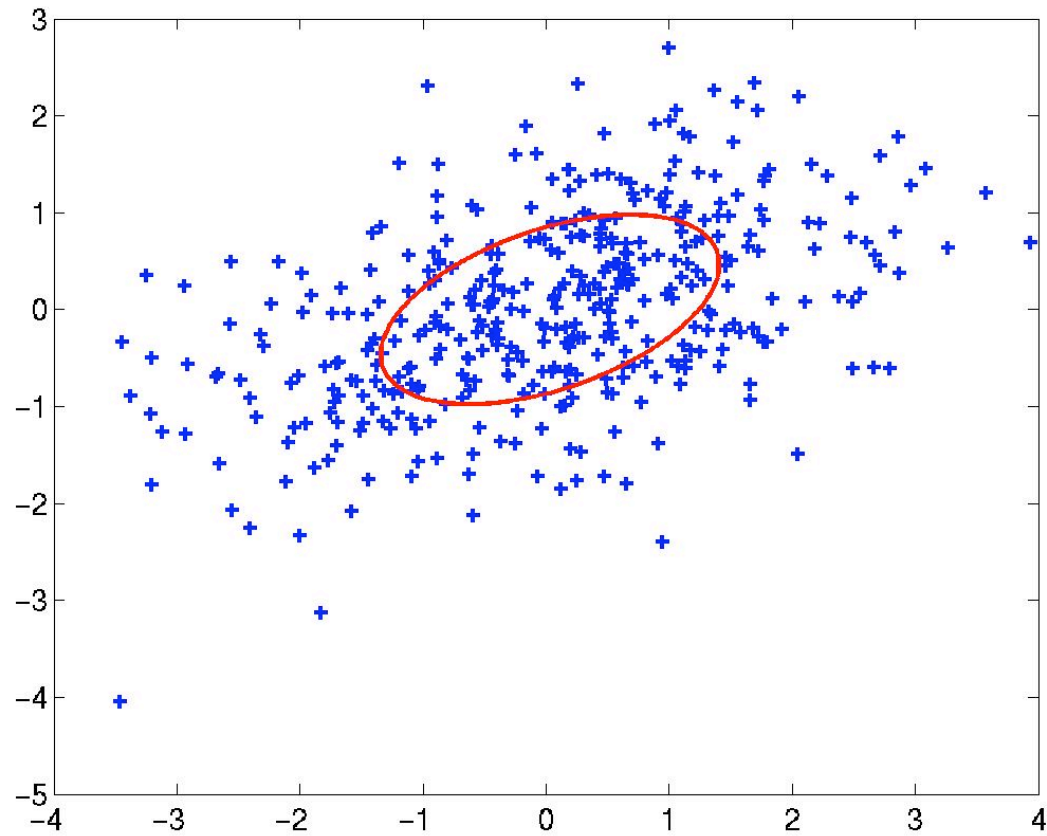
$$p(\mathbf{x} | \Theta) = \frac{1}{2\pi^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)}$$

Here μ is the mean vector and Σ is the covariance matrix

$$\mu = \{\mu_1, \mu_2\} \quad \Sigma = \begin{array}{|c|c|} \hline \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \hline \text{cov}(x_1, x_2) & \text{var}(x_2) \\ \hline \end{array}$$

- The covariance matrix captures linear dependencies among the variables

Example



Important points

- Maximum likelihood estimations (MLE)
- Pseudo counts
- Types of distributions
- Handling continuous variables