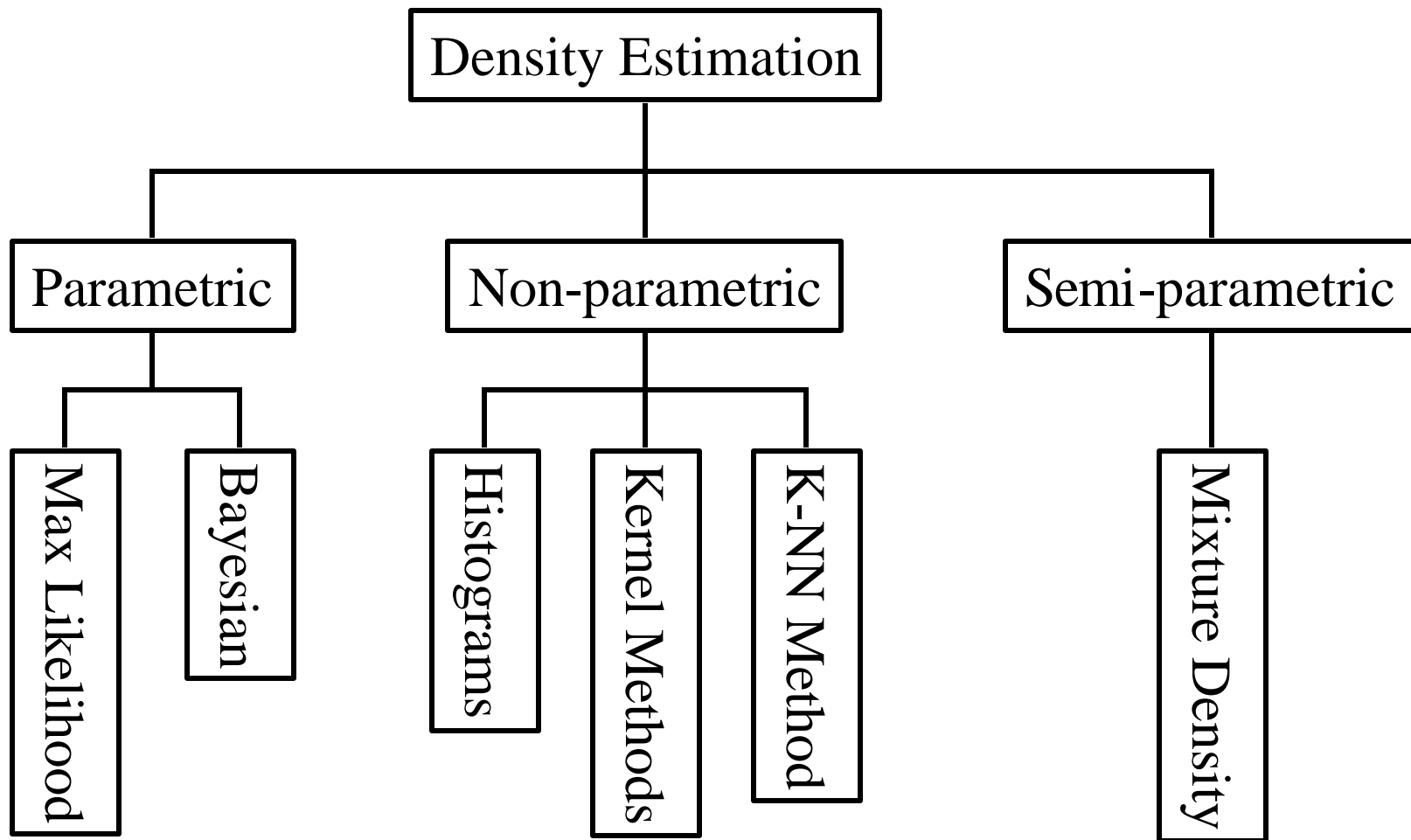


9.913 Pattern Recognition for Vision

Class VI – Density Estimation

Yuri Ivanov

Road Map



Generative vs. Discriminative

There are two schools of thought in Machine Learning:

1. Generative:

- Estimate class models from data
- Compute the discriminant function
- Plug in your data – get the answer

2. Discriminative:

- Estimate the discriminant from data
- Plug in your data – get the answer

Last class

Density Estimation

Density Estimation is at the core of *generative* Pattern Recognition

$$P(a < x < b) = \int_a^b p(x) dx$$

$$\text{mean : } E[x] = \int xp(x) dx$$

$$\begin{aligned} \text{covariance : } E[(x - E[x])(x - E[x])^T] \\ = \int [(x - E[x])(x - E[x])^T] p(x) dx \end{aligned}$$

$$\text{function mean : } E[f(x)] = \int f(x)p(x) dx$$

$$\text{conditional mean : } E[y | x] = \int yp(y | x) dx$$

Refresher

Minimum expected risk:

$$R^* = \int \min_w [R(\mathbf{a} | x)] p(x) dx$$

... is based on conditional risk:

$$\mathbf{w}_i = \arg \min_w R(\mathbf{a} | x)$$

... which is computed from the posterior:

$$R(\mathbf{a} | x) = L(\mathbf{a} | \mathbf{w}) P(\mathbf{w} | x)$$

... which depends on the likelihood:

$$P(\mathbf{w} | x) = \frac{p(x | \mathbf{w}) P(\mathbf{w})}{p(x)}$$

Setting

Data:

$$D = \{D_i\}_{i=1}^C$$

Assume that D_j contains no information about \mathbf{w}_i , $\forall i \neq j$

NOTATIONALLY - we abandon the class label:

$$p(x | \cancel{\mathbf{w}_i}) \implies p(x)$$

Keep in mind:

$$p(x | \mathbf{w}_i) \neq p(x)$$

Goal:

model the probability density function $p(x)$, given a finite number of data points, x_1, x_2, \dots, x_N , drawn from it.

Three Methods

1. Parametric
 - Good: small number of parameters
 - Bad: choice of the parametric form
2. Non-parametric
 - Good: data “dictates” the approximator
 - Bad: large number of parameters
3. Semi-parametric
 - Good: combine the best of both worlds
 - Bad: harder to design
 - Good again: design can be subject to optimization

Parametric Density Estimation

Estimate the density from a given functional family

Given: $p(x | \mathbf{q}) = f(x, \mathbf{q})$

Find: \mathbf{q}

Two methods of parameter estimation:

1. *Maximum Likelihood* method

- Parameters are viewed as unknown but fixed values

2. *Bayesian* method

- Parameters are random variables that have their distributions

Normal (Gaussian) Density Function

A common assumption - *Gaussian*

$$\mathbf{q} = (\mathbf{m}, \Sigma)$$

$$p(x | \mathbf{q}) = \frac{1}{(2\mathbf{p})^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \underbrace{\left((x - \mathbf{m})^T \Sigma^{-1} (x - \mathbf{m})\right)}_{\text{Squared Mahalanobis distance}}\right)$$

Number of
dimensions

“Volume”
of the
covariance

Squared
Mahalanobis
distance

$$\mathbf{m} = E[x]$$

– d parameters

$$\Sigma = E\left[(x - \mathbf{m})^T (x - \mathbf{m})\right]$$

– $d(d + 1) / 2$ parameters

Normal Density

$$\mathbf{m} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

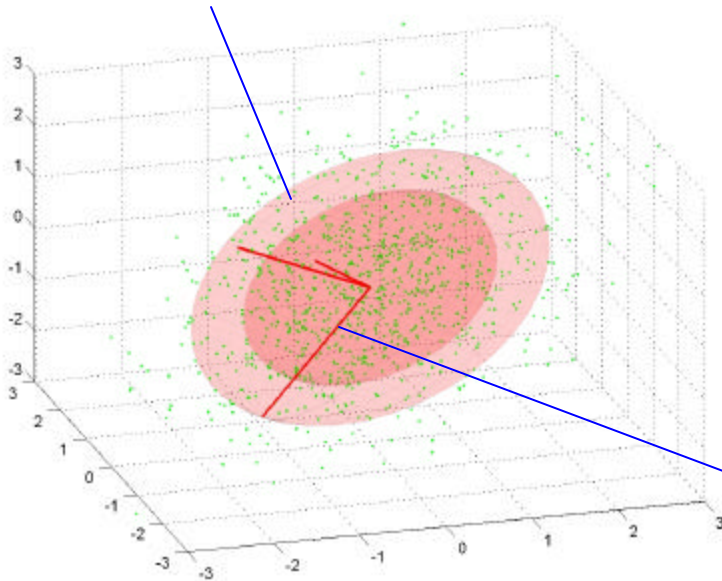
$$\Sigma = \begin{bmatrix} 1 & 0 & .5 \\ 0 & 1 & .3 \\ .5 & .3 & 1 \end{bmatrix}$$

Constant density, $(x - \mathbf{m})^T \Sigma^{-1} (x - \mathbf{m}) = C$ - quadratic surface

Σ - Positive semidefinite



ellipsoid



Principal axes: eigenvectors of Σ
Length: $\sqrt{\mathbf{I}_i}$, \mathbf{I} - eigenvalues of Σ

Whitening Transform

Define:

$$\Lambda = \text{diag}(\text{eigval}(\Sigma)) \quad \text{- Scaling matrix}$$

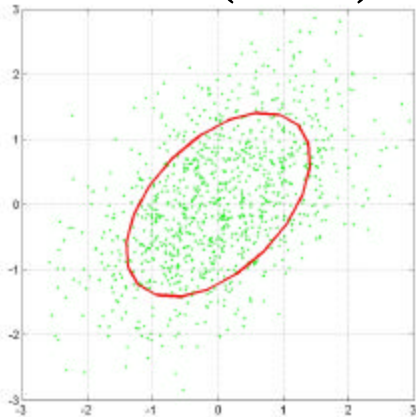
$$\Phi = \text{eigvec}(\Sigma) \quad \text{- Rotation matrix}$$

Then:

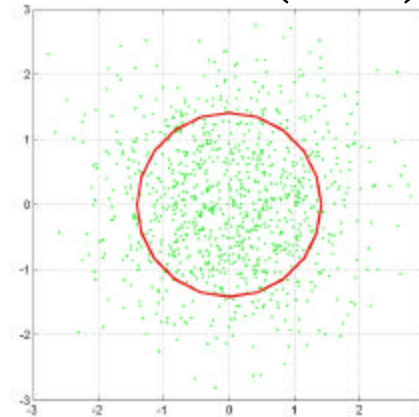
$$W = \Lambda^{-1/2} \Phi^T \quad \text{- "Unscaling" and "unrotates" the data}$$

For all:

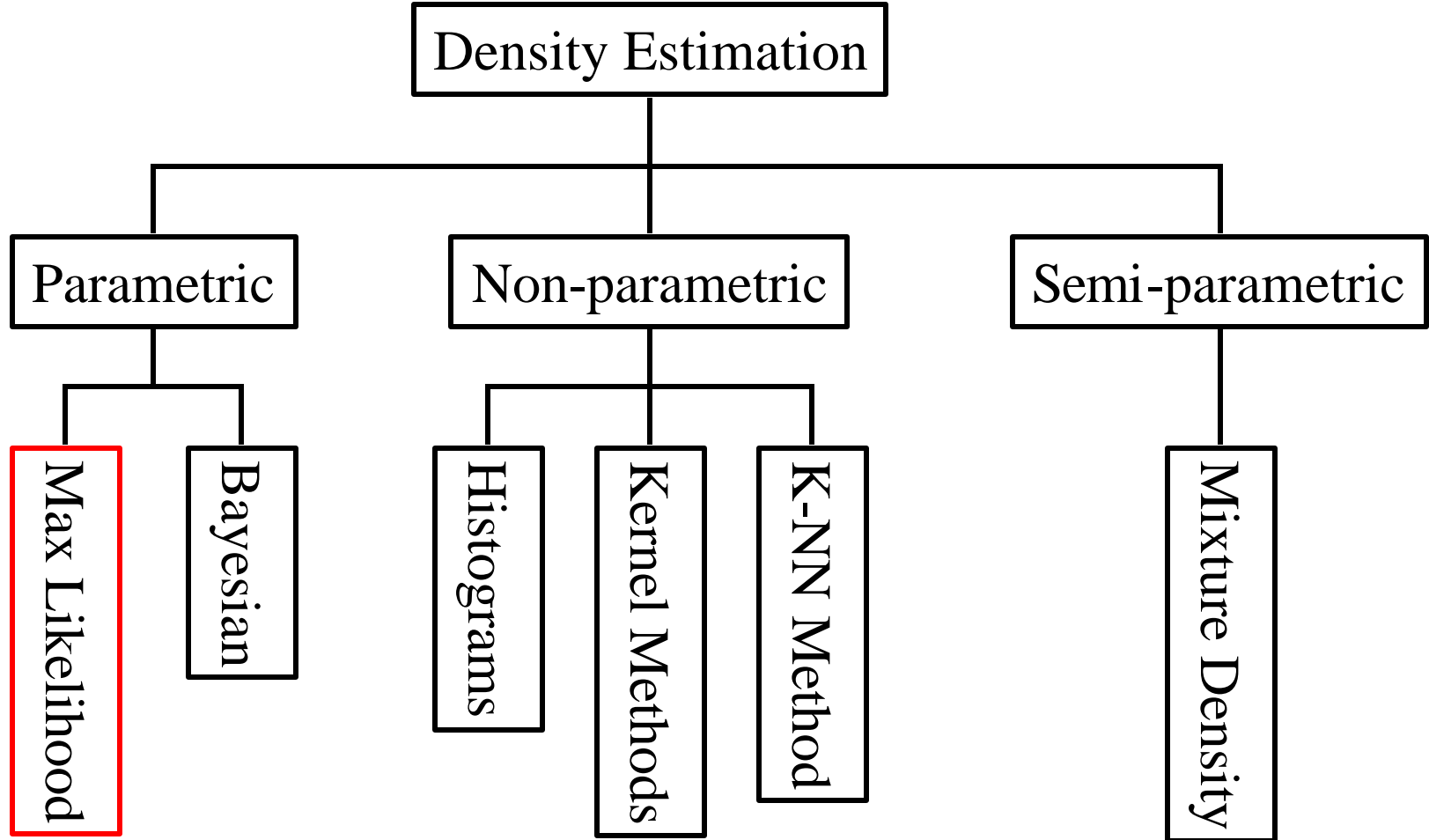
$$x \sim N(\mathbf{0}, \Sigma)$$



$$Wx \sim N(\mathbf{0}, I)$$



You Are Here



Maximum Likelihood

Parameters are fixed but unknown.

$$D \equiv \{x^1, x^2, \dots, x^N\} \quad - \text{ a data set, drawn from } p(x)$$

Notationally, we make density explicitly dependent on parameters:

$$p(x) \implies p(x | \mathbf{q})$$

Assuming that the data is drawn independently (i.i.d.):

$$L(\mathbf{q}) \equiv p(D | \mathbf{q}) = \prod_{n=1}^N p(x^n | \mathbf{q}) \quad - \text{ a likelihood function}$$

To find θ Maximize $L(\theta)$ w.r.t. parameters.

Maximum Likelihood

Maximizing $L(\theta)$ is equivalent to maximizing *log-likelihood function*:

$$l(\mathbf{q}) \equiv \log L(\mathbf{q}) = \log \prod_{n=1}^N p(x^n | \mathbf{q}) = \sum_{n=1}^N \log p(x^n | \mathbf{q})$$

To find θ set the derivative to 0:

$$\nabla_{\mathbf{q}} l(\mathbf{q}) = \sum_{n=1}^N \nabla_{\mathbf{q}} \log p(x^n | \mathbf{q}) = 0$$

And solve for θ

Quick Summary – ML Parameter Estimation

$$P(\mathbf{w}_i | x) = P(\mathbf{w}_i | x, \mathbf{q}_i) = \frac{p(x | \mathbf{w}_i, \mathbf{q}_i) P(\mathbf{w}_i | \mathbf{q}_i)}{p(x | \mathbf{q}_i)} \quad \text{easy}$$

$$p(x | \hat{\mathbf{q}})$$

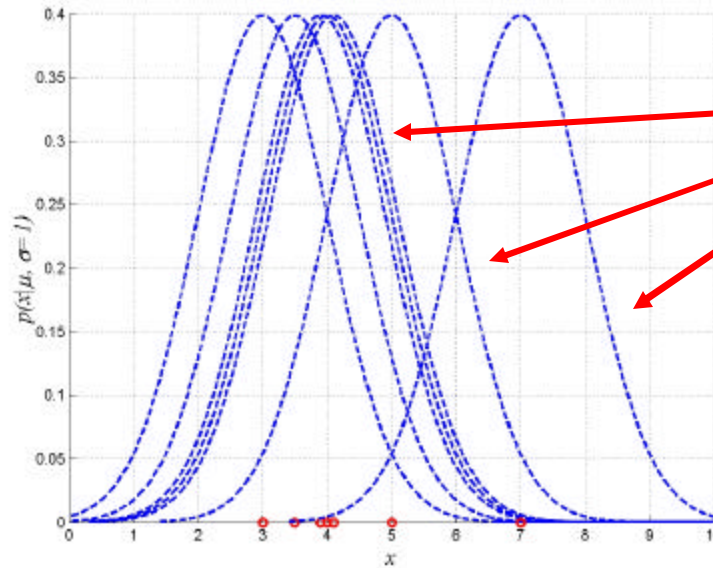
we pick that

$$\operatorname{argmax}_{\mathbf{q}} p(D | \mathbf{q})$$

$$\prod_{n=1}^N p(x^n | \mathbf{q})$$

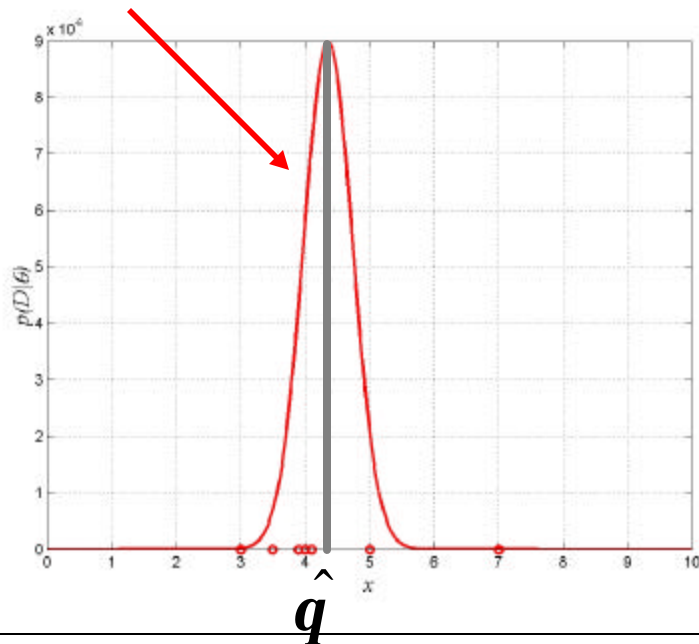
Solving a Maximum Likelihood Problem

Fixed covariance:

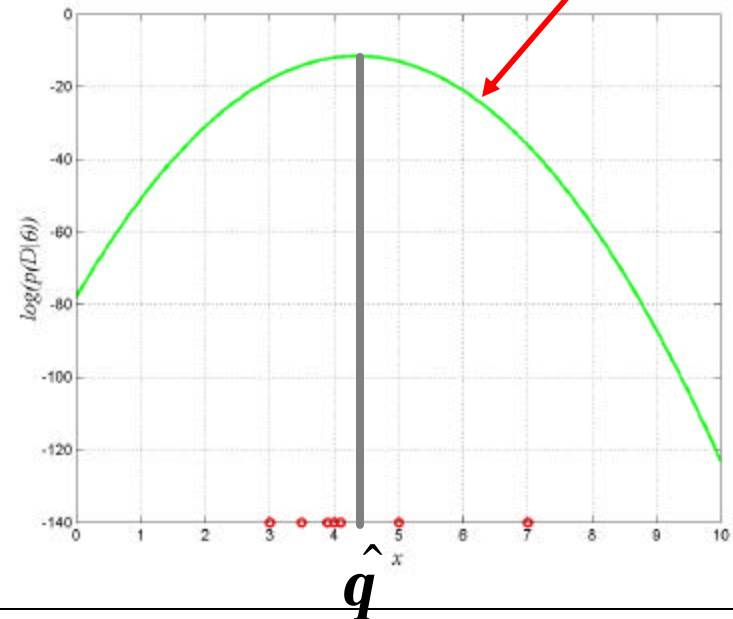


some candidates

$p(D|\mathbf{q})$



$\log p(D|\mathbf{q})$



Maximum Likelihood Example

In d -dimensions:

$$\nabla_{\mathbf{q}} l(\mathbf{q}) = \sum_n \nabla_{\mathbf{q}} \left\{ -\frac{d}{2} \log [2\mathbf{p}] - \frac{1}{2} \log [|\Sigma|] - \frac{1}{2} (x^n - \mathbf{m})^T \Sigma^{-1} (x^n - \mathbf{m}) \right\}$$

Solving for the mean:

$$\nabla_{\mathbf{m}} l(\mathbf{q}) = -\frac{1}{2} \sum_n \Sigma^{-1} (x^n - \hat{\mathbf{m}}) = 0 \Rightarrow$$

$$\hat{\mathbf{m}} = \frac{1}{N} \sum_{n=1}^N x^n$$

- arithmetic average of samples

Maximum Likelihood Example (cont.)

$$\nabla_{\mathbf{q}} l(\mathbf{q}) = \sum_n \nabla_{\mathbf{q}} \left\{ -\frac{d}{2} \log [2\mathbf{p}] - \frac{1}{2} \log [|\Sigma|] - \frac{1}{2} (x^n - \mathbf{m})^T \Sigma^{-1} (x^n - \mathbf{m}) \right\}$$

Solving for the covariance:

For symmetric M : $\frac{d|M|}{dM} = |M| M^{-1}$ and $\frac{d(a^T M^{-1} b)}{dM} = M^{-1} a b^T M^{-1} \implies$

$$\nabla_{\Sigma} l(\mathbf{q}) = -\frac{1}{2} \sum_n \left\{ \hat{\Sigma}^{-1} - \hat{\Sigma}^{-1} (x^n - \hat{\mathbf{m}}) (x^n - \hat{\mathbf{m}})^T \hat{\Sigma}^{-1} \right\} = 0 \implies$$

biased $\hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (x^n - \hat{\mathbf{m}}) (x^n - \hat{\mathbf{m}})^T$

-arithmetic average of
indiv. covariances

Recursive ML

What if data comes one sample at a time?

$$\begin{aligned}\hat{\mathbf{m}}_N &= \frac{1}{N} \sum_{n=1}^N x^n = \frac{1}{N} \left[x^N + \sum_{n=1}^{N-1} x^n \right] \\ &= \frac{1}{N} \left[x^N + (N-1) \hat{\mathbf{m}}_{N-1} \right] = \hat{\mathbf{m}}_{N-1} + \frac{1}{N} \left[x^N - \hat{\mathbf{m}}_{N-1} \right]\end{aligned}$$

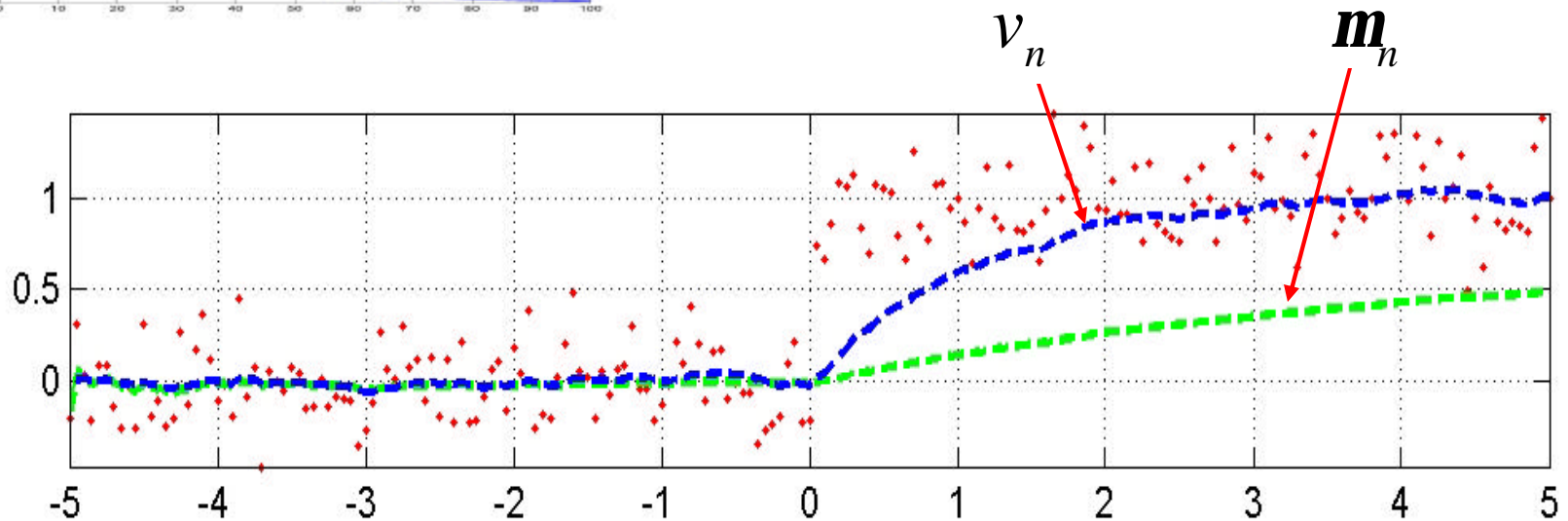
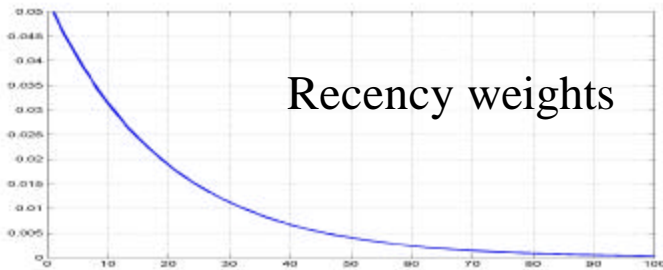
This estimate “stiffens” with more data (as it should).

One idea – fix the fraction. Then the estimate can track a non-stationary process

Recursive ML

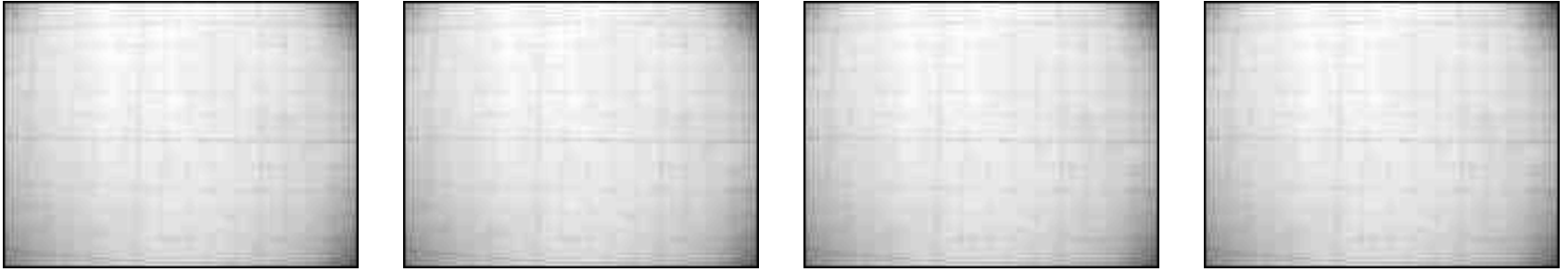
Fix the update rate and retrace the steps:

$$\begin{aligned}v_N &= v_{N-1} + \mathbf{g} \left[x^N - v_{N-1} \right] = (1 - \mathbf{g})v_{N-1} + \mathbf{g}x^N \\ &= (1 - \mathbf{g})^2 v_{N-2} + (1 - \mathbf{g})\mathbf{g}x^{N-1} + \mathbf{g}x^N \\ &= (1 - \mathbf{g})^M v_{N-M} + \sum_{k=1}^M (1 - \mathbf{g})^{M-k} \mathbf{g}x^k\end{aligned}$$



Simple Example

Several images from a static camera:



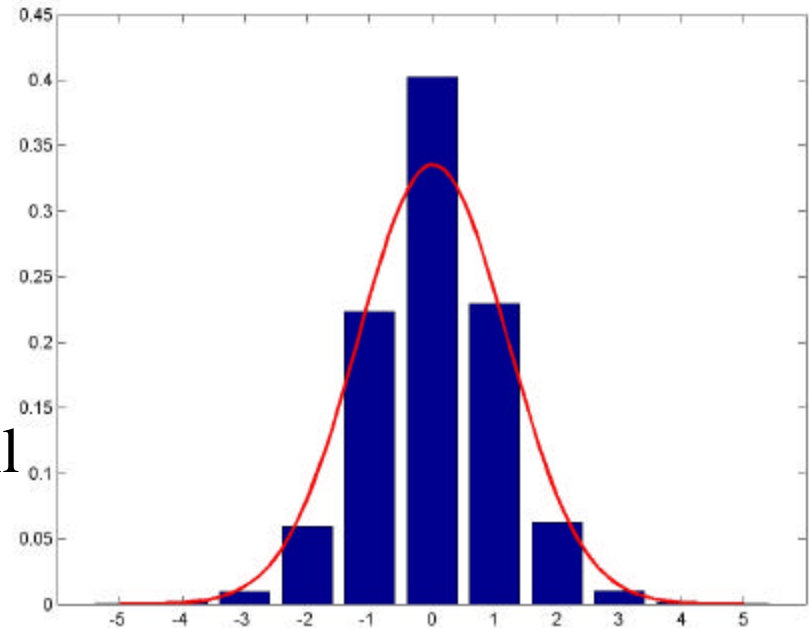
How much noise is in it?

$$x = \text{vec}(I_t - I_{t-1})$$

$$\mathbf{m} = 0$$

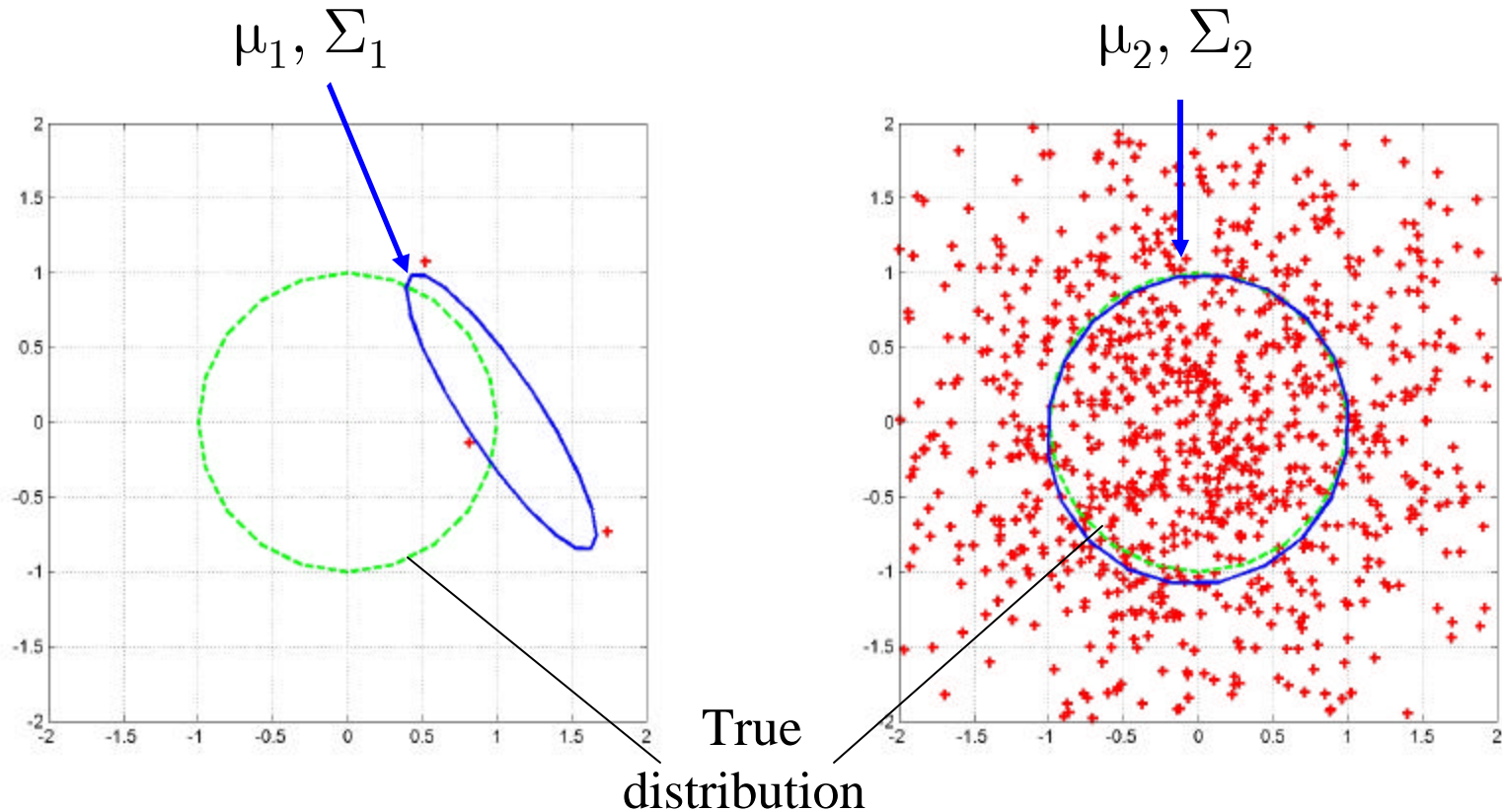
$$\mathbf{s} = 1.2$$

Now we can set a threshold that will *statistically* distinguish pixel noise from an object



Problems with ML

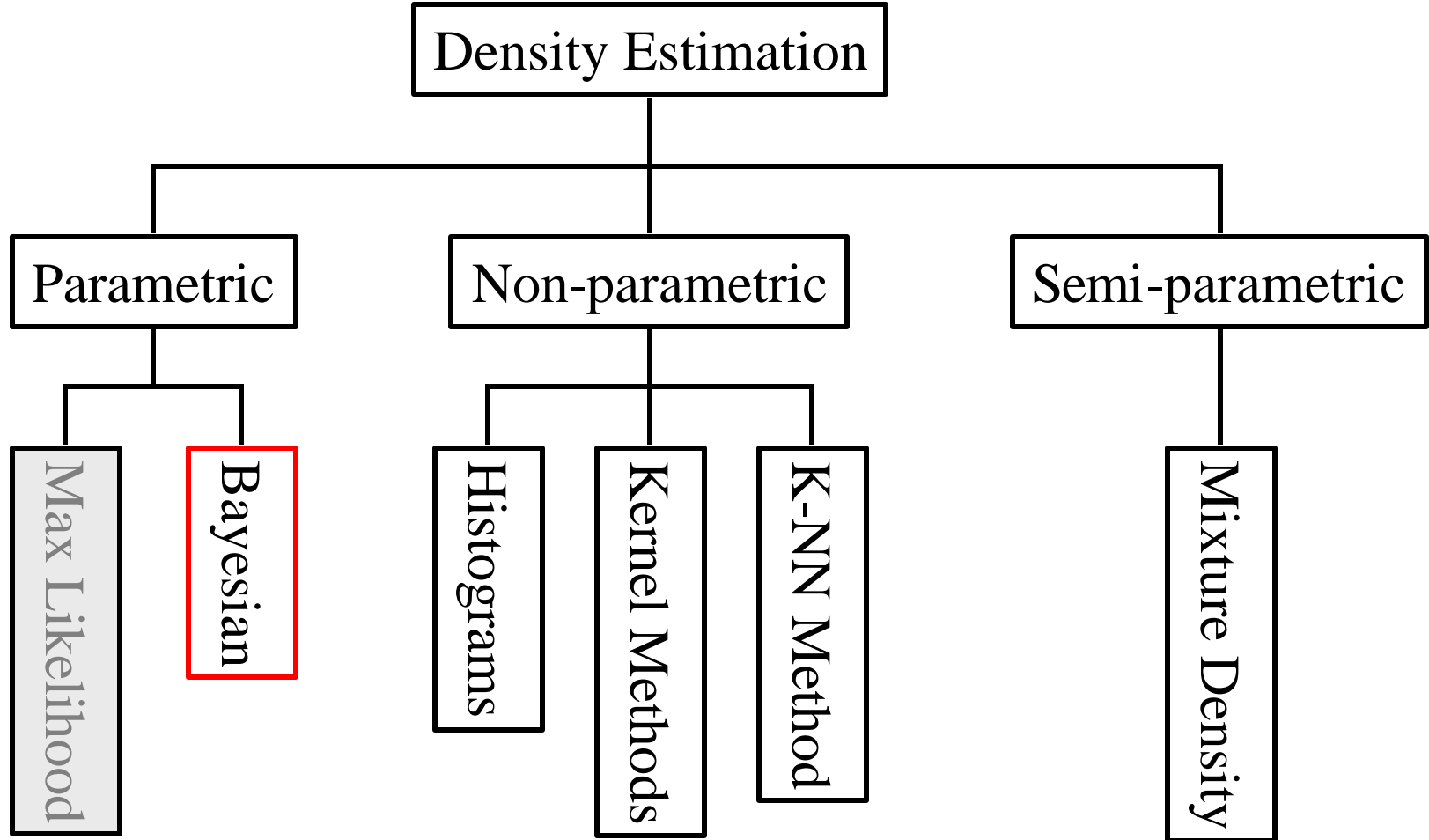
We are given two estimates:



Which one do we believe?

ML gives a single solution, regardless of uncertainty.

You Are Here



Bayesian Parameter Estimation

In classification our goal so far has been to estimate $P(\mathbf{w} | x)$

Let's make the dependency on the data explicit:

$$P(\mathbf{w}_i | x, D) = \frac{p(x | \mathbf{w}_i, D)P(\mathbf{w}_i | D)}{p(x | D)}$$

$P(\mathbf{w}_i | D)$ - this is easy to compute

$P(x | D)$ - this is easy to compute by marginalization

What about $p(x | \mathbf{w}_i, D)$?

Bayesian Parameter Estimation



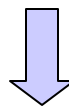
$$P(\mathbf{w}_i | x, D) = \frac{p(x | \mathbf{w}_i, D)P(\mathbf{w}_i | D)}{p(x | D)}$$

This is a supervised problem so far:

$$D = \{D_1, D_2, \dots, D_N\}$$

$$p(x | \mathbf{w}_i, D) = p\left(x | \mathbf{w}_i, \{D_j\}_{j=1 \dots N}\right)$$

$$= p\left(x | \mathbf{w}_i, D_i, \{\cancel{D_j}\}_{j \neq i}\right) = p(x | \mathbf{w}_i, D_i)$$



$$P(\mathbf{w}_i | x, D) = \frac{p(x | \mathbf{w}_i, D_i)P(\mathbf{w}_i | D)}{p(x | D)}$$

Bayesian Parameter Estimation

We will assume that we can obtain “labeled” data, so again:

Notationally: $p(x | \mathbf{w}_i, D_i) \implies p(x | D)$

Now our problem is to compute density for x given the data D .

We assume the form of $p(x)$ – the source density for D :

$$p(x) \implies p(x | \mathbf{q})$$

... and treat θ as a *random variable*

Bayesian Parameter Estimation

Instead of choosing a value for a parameter, we use them all:

$$p(x | D) = \int p(x, \mathbf{q} | D) d\mathbf{q} = \int p(x | \mathbf{q}, \cancel{D}) p(\mathbf{q} | D) d\mathbf{q}$$

Data predicts the new sample

x is independent of D given \mathbf{q}

$$= \int \underbrace{p(x | \mathbf{q})}_{\text{We chose the form of this}} \underbrace{p(\mathbf{q} | D)}_{\text{What is this?}} d\mathbf{q}$$

*We chose the form
of this*

What is this?

Average densities $p(x/\mathbf{q})$ for ALL possible values of \mathbf{q} weighted by its posterior probability

Bayesian Parameter Estimation



$$\int p(x | \mathbf{q}) p(\mathbf{q} | D) d\mathbf{q}$$

Computing the posterior probability for \mathbf{q} :

Using Bayes rule:

What is this?

*Prior belief about
the parameters
(density)*

$$p(\mathbf{q} | D) = \frac{p(D | \mathbf{q}) p(\mathbf{q})}{\int p(D | \mathbf{q}) p(\mathbf{q}) d\mathbf{q}}$$

Using independence:

$$p(D | \mathbf{q}) = \prod_{n=1}^N p(x^n | \mathbf{q})$$

Bayesian method does not commit to a particular value of θ , but uses the entire distribution.

Quick Summary – Bayesian Parameter Estimation

$$P(\mathbf{w}_i | x) = P(\mathbf{w}_i | x, D) = \frac{p(x | \mathbf{w}_i, D_i) P(\mathbf{w}_i | D)}{p(x | D)} \quad \text{easy}$$

hard

$$\int p(x | \mathbf{q}) p(\mathbf{q} | D) d\mathbf{q}$$

we pick that

$$\frac{p(D | \mathbf{q}) p(\mathbf{q})}{p(D)}$$

we “know” this*,**

hard

$$\prod_{n=1}^N p(x^n | \mathbf{q})$$

* *Non-informative* prior – doesn’t introduce bias

** *Conjugate prior* – causes $p(\mathbf{q} | D)$ have the same functional form as $p(D | \mathbf{q})$

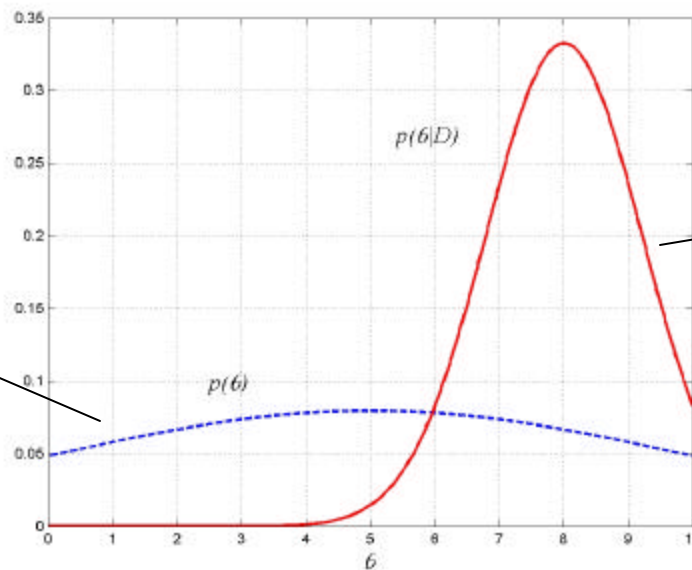
Bayesian Parameter Estimation



$$\int p(x | \mathbf{m}) p(\mathbf{m} | D) d\mathbf{m}$$

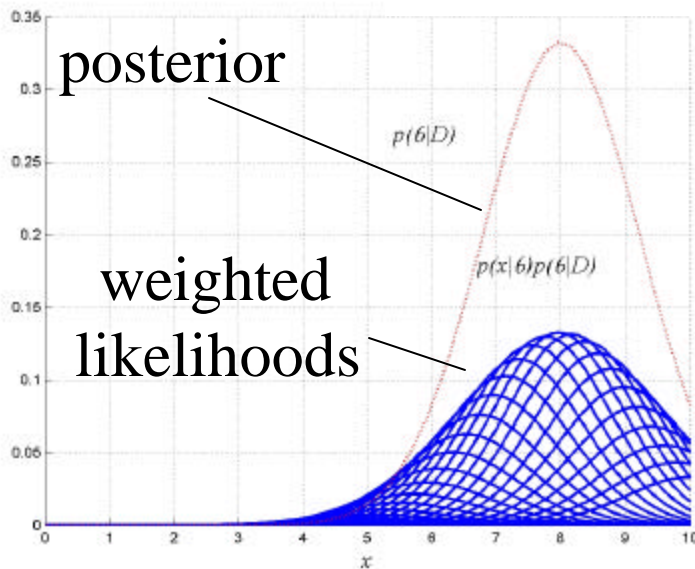
For $\mathbf{q} = \mathbf{m}$:

Parameter prior



Parameter posterior

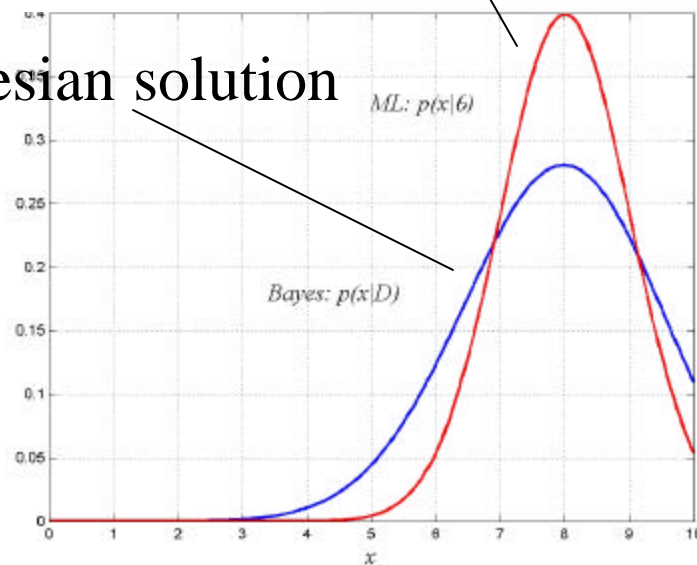
ML solution



posterior

weighted likelihoods

Bayesian solution



Bayesian Parameter Estimation - Example



First let's deal with the parameter:

$$\int p(x | \mathbf{m}) p(\mathbf{m} | D) d\mathbf{m}$$

Likelihood: $p(x | \mathbf{m}) = \mathbf{N}(\mathbf{m}, \mathbf{s}^2)$ *fixed*

Parameter prior: $p(\mathbf{m}) = \mathbf{N}(\mathbf{m}_0, \mathbf{s}_0^2)$

Need to find: $p(\mathbf{m} | D)$

Bayes rule again:

$$p_N(\mathbf{m} | D) = \frac{p(D | \mathbf{m}) p(\mathbf{m})}{p(D)} = \underbrace{\mathbf{a} \left[\prod_{n=1}^N p(x^n | \mathbf{m}) \right] \mathbf{N}(\mathbf{m}_0, \mathbf{s}_0^2)}_{\text{This is a Gaussian}} = \mathbf{N}(\mathbf{m}_N, \mathbf{s}_N)$$

N-sample parameter posterior

This is a Gaussian

Need these

Bayesian Parameter Estimation - Example

So, the posterior is a Gaussian

$$p_N(\mathbf{m} | D) = \mathbf{N}(\mathbf{m}_N, \mathbf{s}_N)$$

After some algebra and identifying the terms:

$$\frac{1}{\mathbf{s}_N^2} = \frac{1}{\mathbf{s}^2} N + \frac{1}{\mathbf{s}_0^2} \quad - \text{when Gaussians multiply – precisions add}$$

... and

$$\mathbf{m}_N = \frac{N \mathbf{s}_0^2}{N \mathbf{s}_0^2 + \mathbf{s}^2} \bar{x} + \frac{\mathbf{s}^2}{N \mathbf{s}_0^2 + \mathbf{s}^2} \mathbf{m}_0$$

With increasing N covariance of the posterior decreases and the prior becomes unimportant.

Bayesian Parameter Estimation - Example

Now the integral:

$$p(x | D) = \int p(x | \mathbf{q}) p(\mathbf{q} | D) d\mathbf{q}$$
$$= \int \mathbf{N}(\mathbf{m}, \mathbf{S}^2) \mathbf{N}(\mathbf{m}_N, \mathbf{S}_N^2) d\mathbf{m} = \mathbf{N}(\mathbf{m}_N, \mathbf{S}^2 + \mathbf{S}_N^2)$$

↓
*You can show that it
is also a Gaussian*

Any guesses about why Gaussian is such a common assumption?

Recursive Bayes



$$\int p(x | \mathbf{q}) p(\mathbf{q} | D) d\mathbf{q}$$

For N -point likelihood:

$$\begin{aligned} p(D^N | \mathbf{q}) &= \prod_{n=1}^N p(x^n | \mathbf{q}) \\ &= p(x^N | \mathbf{q}) \prod_{n=1}^{N-1} p(x^n | \mathbf{q}) = p(x^N | \mathbf{q}) p(D^{N-1} | \mathbf{q}) \end{aligned}$$

From this the recursive relation for the posterior:

$$\begin{aligned} p(\mathbf{q} | D^N) &= \frac{p(x^N | \mathbf{q}) p(D^{N-1} | \mathbf{q}) p(\mathbf{q})}{p(D^N)} \\ &= \frac{p(x^N | \mathbf{q}) p(\mathbf{q} | D^{N-1})}{\int p(x^N | \mathbf{q}) p(\mathbf{q} | D^{N-1}) d\mathbf{q}} \end{aligned}$$

Recursive Bayes (cont.)

Again:

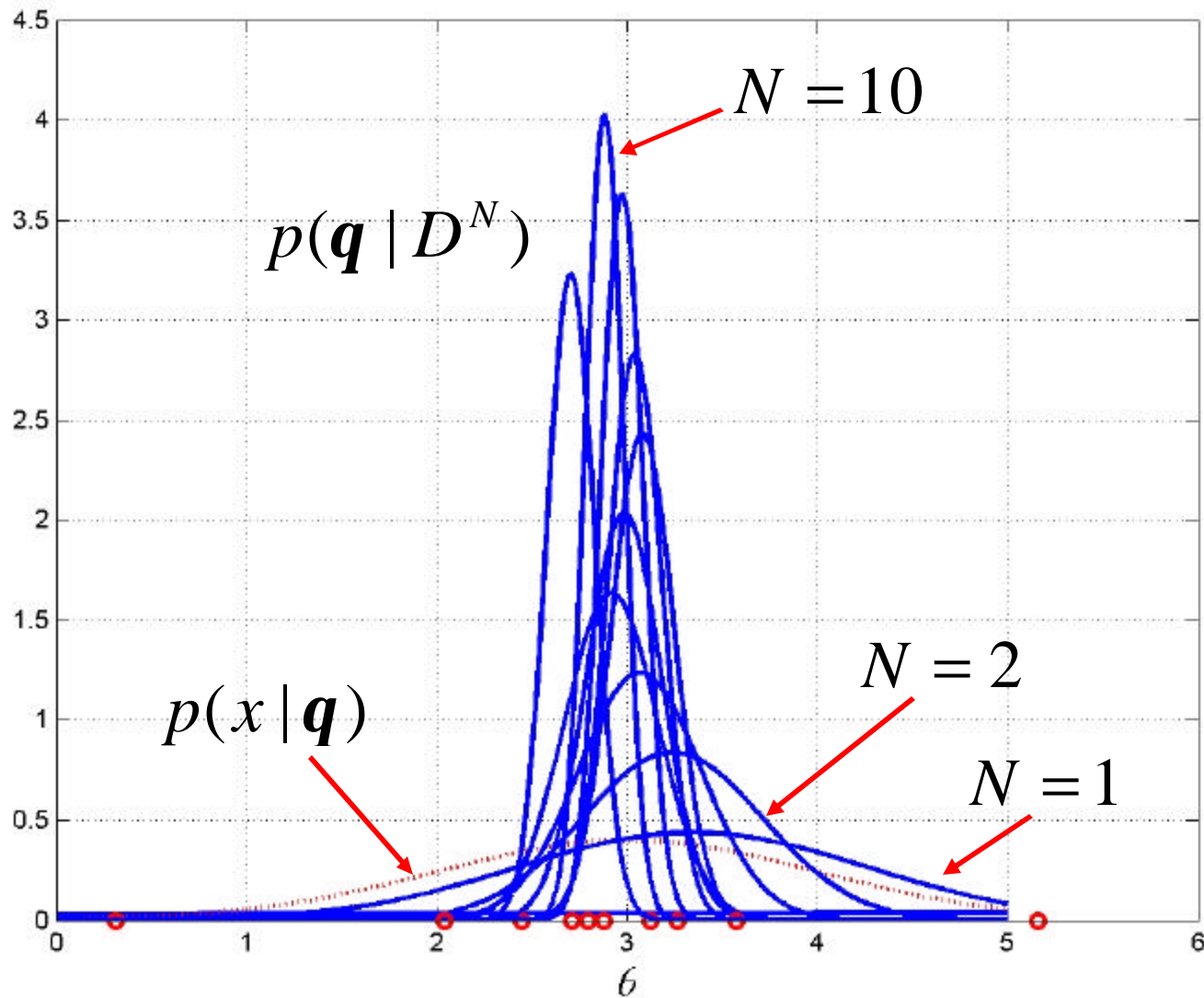
$$p(\mathbf{q} | D^N) = \frac{p(x^N | \mathbf{q}) p(\mathbf{q} | D^{N-1})}{\int p(x^N | \mathbf{q}) p(\mathbf{q} | D^{N-1}) d\mathbf{q}} \quad - \text{1-point update.}$$

Setting $N=1$:

$$\frac{1}{\mathbf{s}_n^2} = \frac{1}{\mathbf{s}^2} + \frac{1}{\mathbf{s}_{n-1}^2}$$

$$\mathbf{m}_n = \frac{\mathbf{s}_{n-1}^2}{\mathbf{s}_{n-1}^2 + \mathbf{s}^2} x + \frac{\mathbf{s}^2}{\mathbf{s}_{n-1}^2 + \mathbf{s}^2} \mathbf{m}_{n-1}$$

Recursive Bayes (cont.)



Problems with Bayesian Method

1. Integration is difficult
2. Analytic solutions are only available for restricted class of densities
3. Technicality: If the true $p(x/\mathbf{q})$ is NOT what we assume it is, the prior probability of any parameter setting is 0!
4. Integration is difficult
5. Did I mention that the integration is hard?

Relation between Bayesian and ML Inference

$$p(\mathbf{q} | D) \propto p(D | \mathbf{q}) p(\mathbf{q})$$

peaks at $\hat{\mathbf{q}}_{ML}$

$$= \left[\prod_n p(x^n | \mathbf{q}) \right] p(\mathbf{q}) = L(\mathbf{q}) p(\mathbf{q})$$

If the peak is sharp and $p(\theta)$ is flat, then:

$$p(x | D) = \int p(x | \mathbf{q}) p(\mathbf{q} | D) d\mathbf{q}$$
$$\simeq \int p(x | \hat{\mathbf{q}}) p(\mathbf{q} | D) d\mathbf{q} = p(x | \hat{\mathbf{q}}) \int p(\mathbf{q} | D) d\mathbf{q} = p(x | \hat{\mathbf{q}})$$

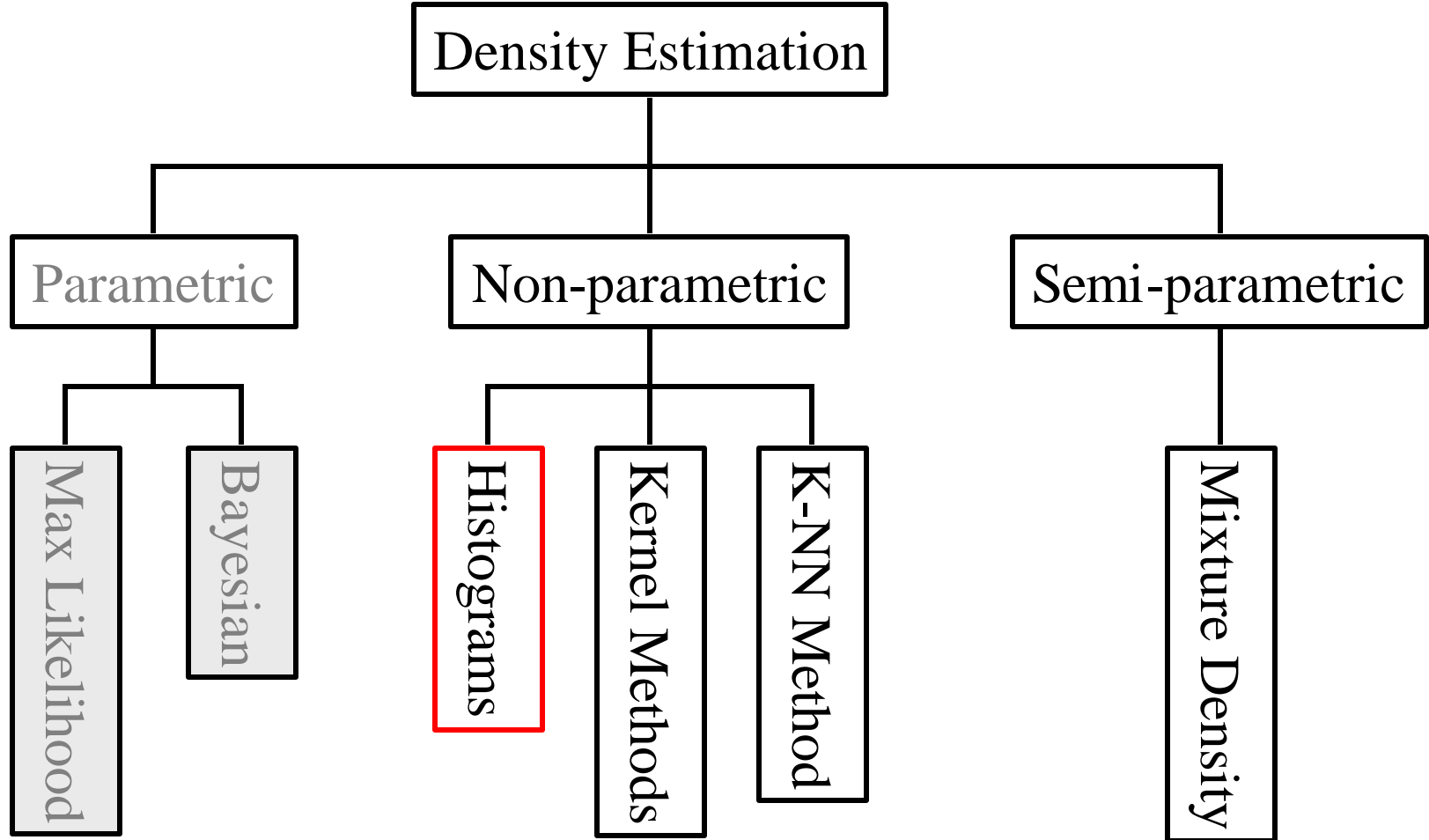
$$\text{As } N \rightarrow \infty, p(x | D) \leftrightarrow p(x | \hat{\mathbf{q}})$$

Non-Parametric Methods for Density Estimation

Non-parametric methods do not assume any particular form for $p(x)$

1. Histograms
2. Kernel Methods
3. K-NN method

You Are Here



Histograms

$\hat{P}(x)$ is a discrete approximation of $p(x)$

- Count a number of times that x lands in the i -th bin

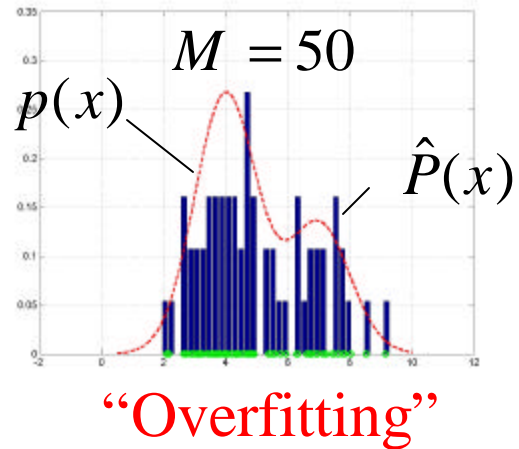
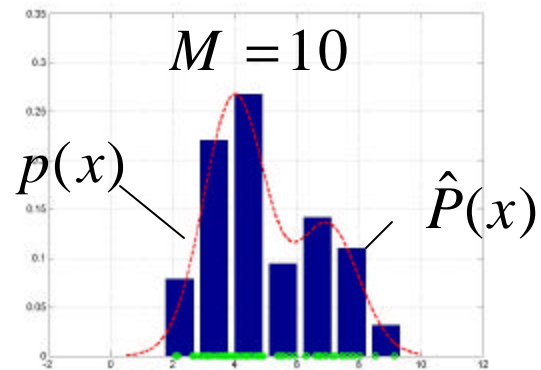
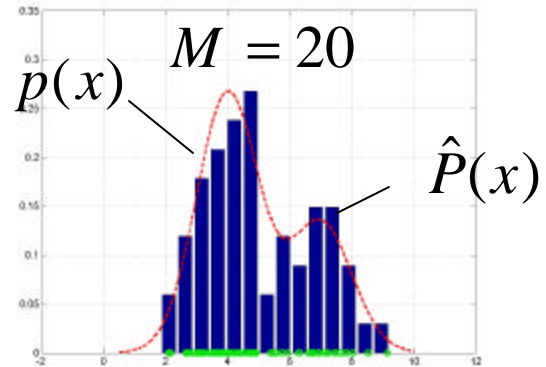
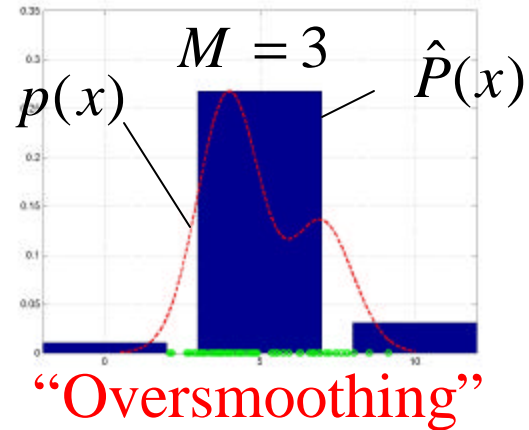
$$H(i) = \sum_{j=1}^N I(x \in R_i), \quad \forall i = 1, 2, \dots, M$$

- Normalize

$$\hat{P}(i) = \frac{H(i)}{\sum_{j=1}^M H(j)}$$

Histograms

How many bins?



Histograms

Good:

- Once it is constructed, the data can be discarded
- Quick and intuitive

Bad:

- Very sensitive to number of bins, M
- Estimated density is not smooth
- Poor generalization in higher dimensions

Aside: Curse of dimensionality (Bellman, '61):

- Imagine we build a histogram of a 1-d feature (say, *Hue*)
 - 10 bins
 - 1 bin = 10% of the input space
 - need at least 10 points to populate every bin
- We add another feature (say, *Saturation*)
 - 10 bins again
 - 1 bin = 1% of the input space
 - we need at least 100 points to populate every bin
- We add another feature (say, *Value*)
 - 10 bins again
 - 1 bin = 0.1% of the input space
 - we need at least 1000 points to populate every bin

$$N = b^d$$

- number of points grows exponentially

Aside: Curse Continues

Volume of a cube in R^d with side l :

$$V_l = l^d$$

Volume of a cube with side $l-\epsilon$:

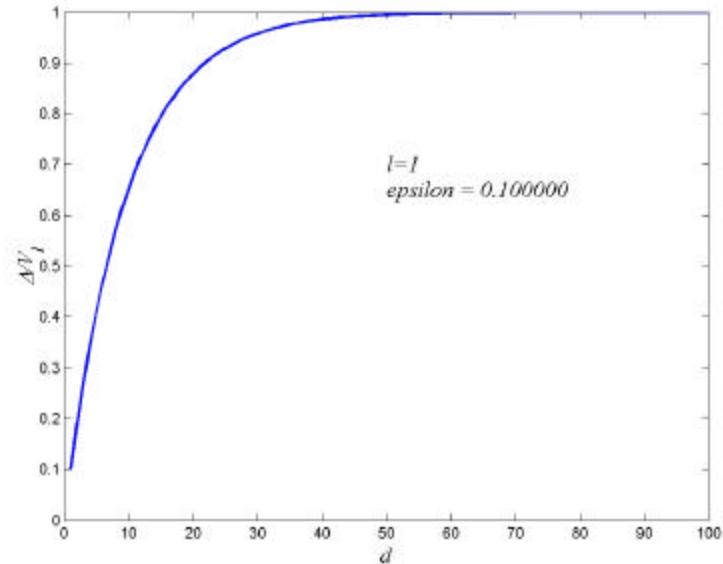
$$V_\epsilon = (l - \epsilon)^d$$

Volume of the ϵ -shell:

$$\Delta = V_l - V_\epsilon = l^d - (l - \epsilon)^d$$

Ratio of the volume of the ϵ -shell to the volume of the cube:

$$\frac{\Delta}{V_l} = \frac{l^d - (l - \epsilon)^d}{l^d} = 1 - \left(1 - \frac{\epsilon}{l}\right)^d \rightarrow 1 \text{ as } d \rightarrow \infty \text{ !!!!!}$$



Aside: Lessons of the curse

In generative models:

- Use as much data as you can get your hands on
- Reduce dimensionality as much as you can get away with

<End of Digression>

General Reasoning

By definition:

$$P(x \in R) = P = \int_R p(x') dx'$$

If we have N i.i.d. points drawn from $p(x)$:

$$P(|x \in R| = k) = \frac{N!}{k!(N-k)!} P^k (1-P)^{N-k} = B(N, P)$$

Num. of unique splits
K vs. (N-K)

Prob that k of
particular x -es are in R

Prob that the
rest are not

$B(N, P)$ is a *binomial* distribution of k

General Reasoning (cont.)

Mean and variance of $B(N, P)$:

$$\text{Mean: } \mathbf{m} = E[k] = NP \quad \Rightarrow P = E[k / N]$$

$$\text{Variance: } \mathbf{s}^2 = E[(k - \mathbf{m})^2] = NP(1 - P)$$

$$\Rightarrow E\left[(k / N - P)^2\right] = \left(\frac{\mathbf{s}}{N}\right)^2 = P(1 - P) / N$$

That is:

- $E[k/N]$ is a good estimate of P
- P is distributed around this estimate with vanishing variance

So:

$$P \simeq k / N$$

General Reasoning (cont.)

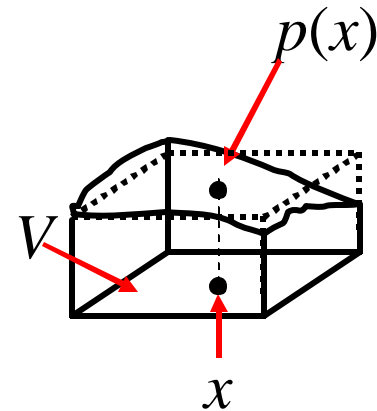
So:

$$P \simeq k / N$$

On the other hand, under mild assumptions:

$$P = \int_R p(x') dx' \simeq p(x)V$$

Volume of R
(not $p(x)$)



... which leads to:

$$p(x) \simeq \frac{k}{NV}$$

General Reasoning (cont.)

Now, given N data points – how do we really estimate $p(x)$?

$$p(x) \approx \frac{k}{NV}$$

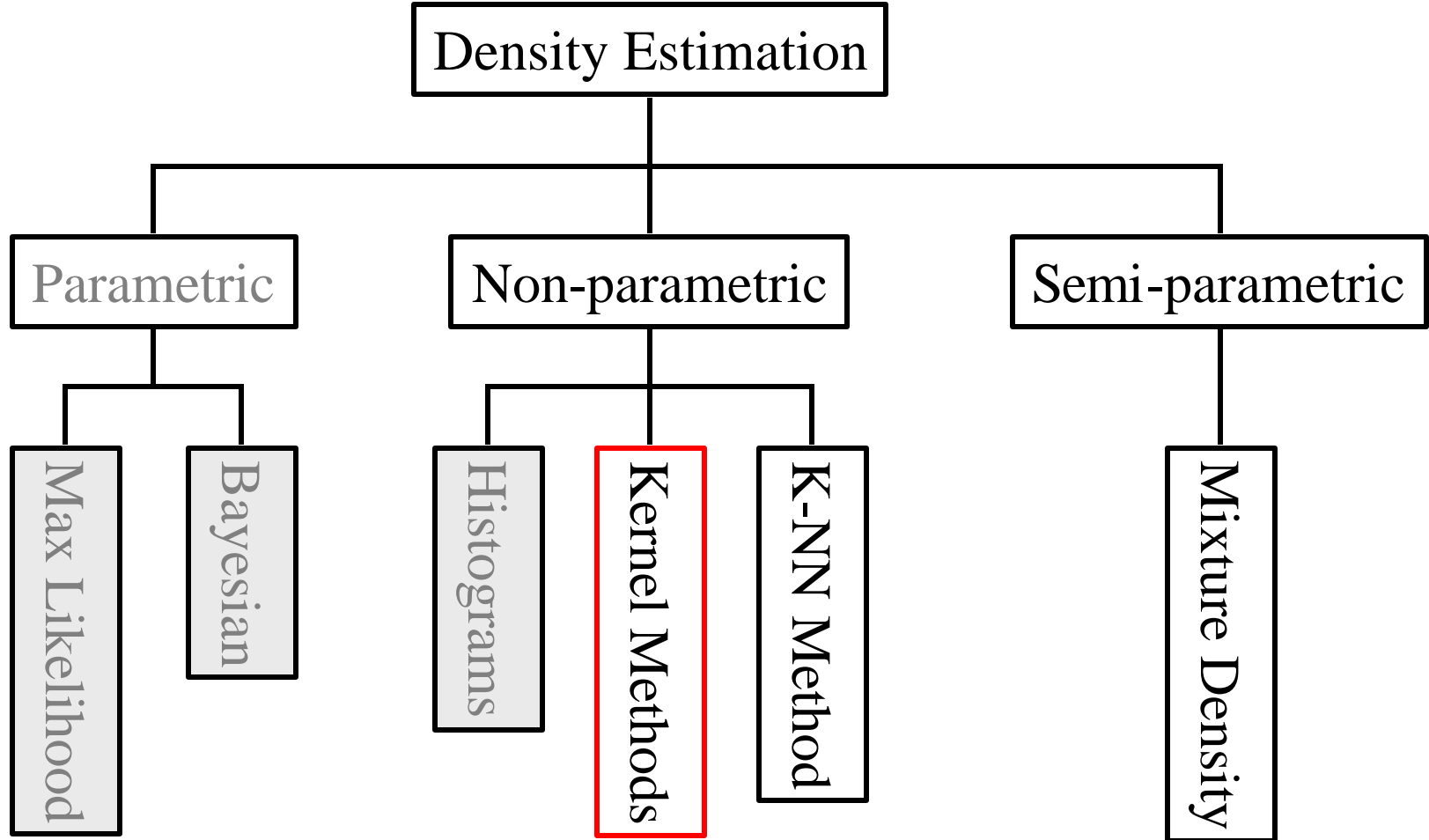
Fix k and vary V
until it encloses k
points

K-Nearest Neighbors (KNN)

Fix V and count how
many points (k) it
encloses

Kernel methods

You Are Here



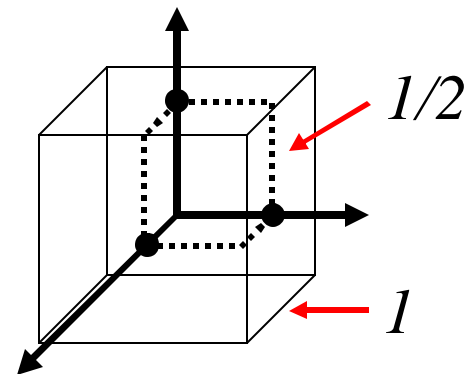
Kernel Methods of Density Estimation

We choose V by specifying a hypercube with a side h :

$$V = h^d$$

Mathematically:

$$H(\mathbf{y}) = \begin{cases} 1 & |y_j| < 1/2 \quad j = 1, \dots, d \\ 0 & \text{otherwise} \end{cases}$$



kernel function:

$$H(\mathbf{y}) \geq 0, \quad \forall \mathbf{y} \quad \text{and} \quad \int H(\mathbf{y}) d\mathbf{y} = 1$$

Parzen Windows

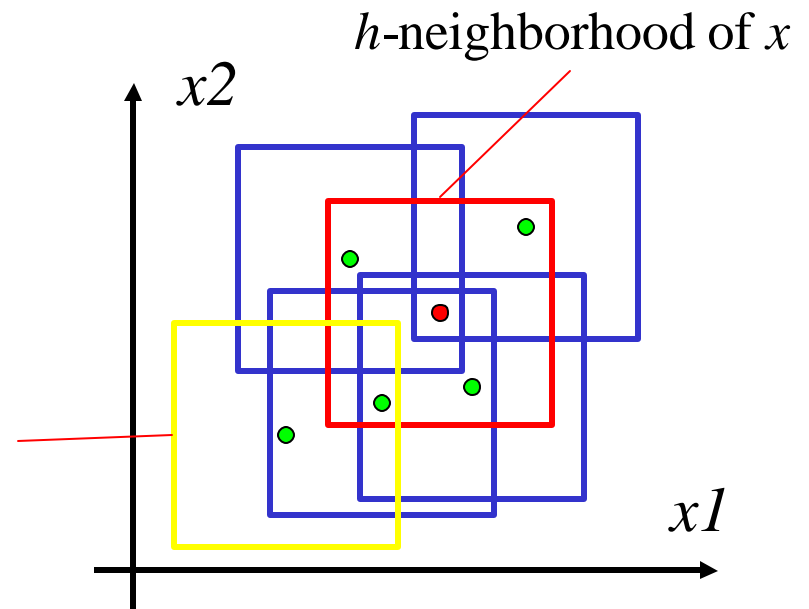
Then

$H\left(\left(\mathbf{x} - \mathbf{x}^n\right) / h\right)$ - a hypercube with side h centered at \mathbf{x}^n

H can help count the points in a volume V around any x :

$$k(x) = \sum_{n=1}^N H\left(\frac{x - x^n}{h}\right)$$

No contribution
to the count at x



Rectangular Kernel

So the number of points in h -neighborhood of x :

$$k(x) = \sum_{n=1}^N H\left(\frac{x - x^n}{h}\right)$$

... is easily converted to the density estimate:

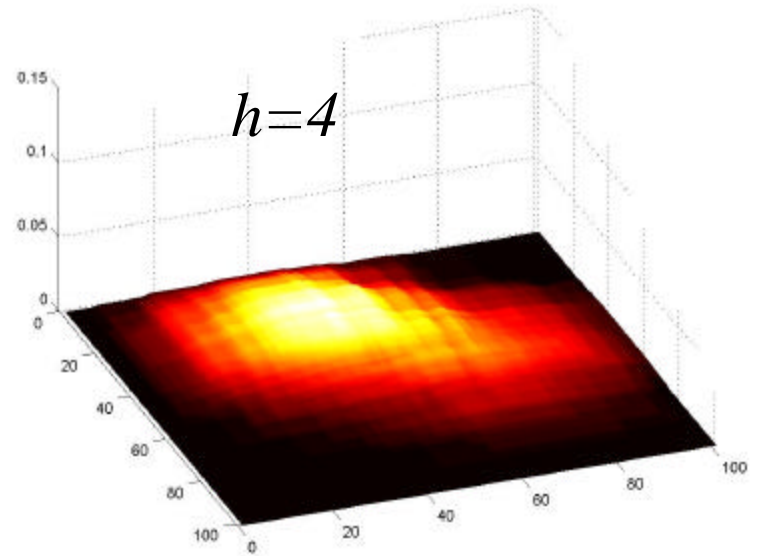
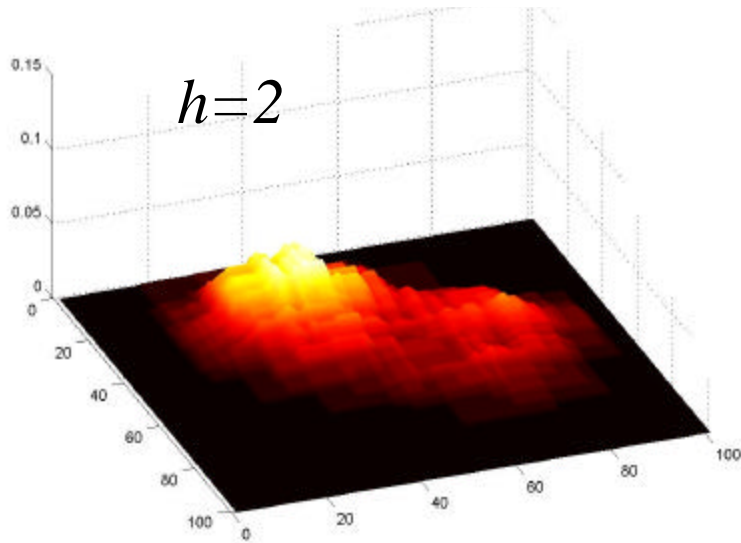
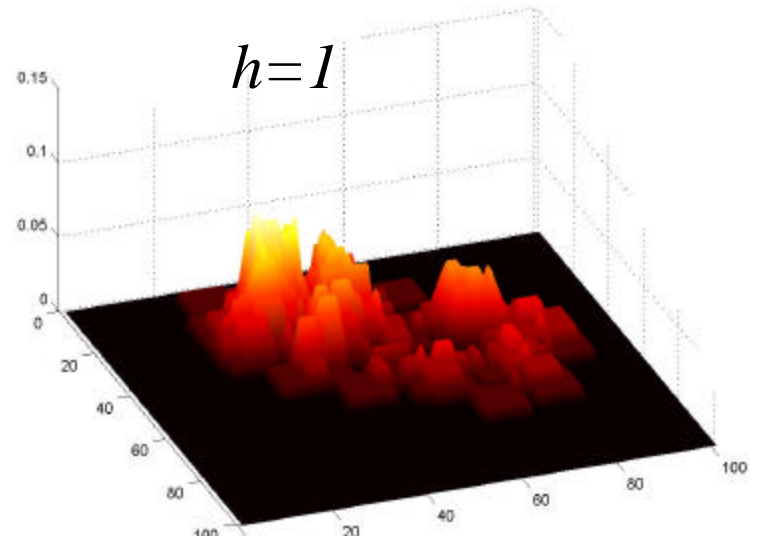
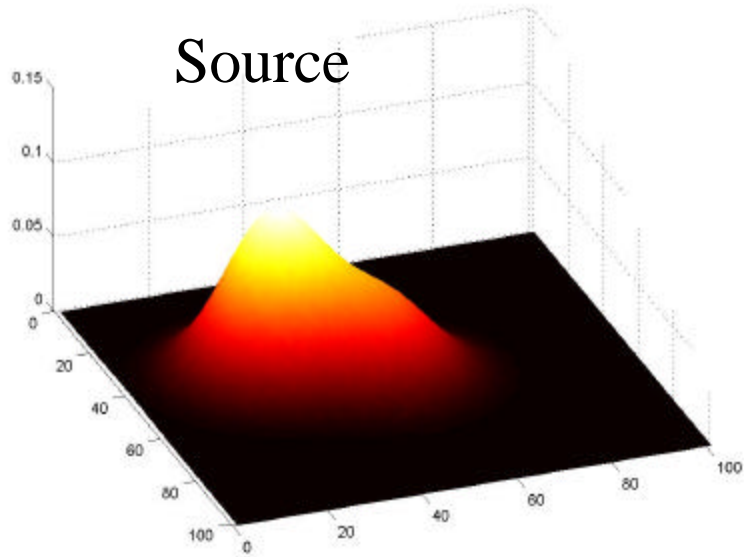
$$\tilde{p}(x) = \frac{k(x)}{NV} = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^d} H\left(\frac{x - x^n}{h}\right) \leftarrow K(x, x^n)$$

Integrates to 1

Subtle point:

$$\int \left[\frac{1}{N} \sum_{n=1}^N K(x, x^n) \right] dx = \frac{1}{N} \sum_{n=1}^N \left[\int K(x, x^n) dx \right] = 1$$
$$\Rightarrow \int \tilde{p}(x) dx = 1$$

Example



Smoothed Window Functions

The problem is as in histograms – it is discontinuous

We can choose a smoother function, s.t.:

$$\tilde{p}(x) \geq 0, \quad \forall x \quad \text{and} \quad \int \tilde{p}(x) dx = 1$$

Ensured by kernel conditions

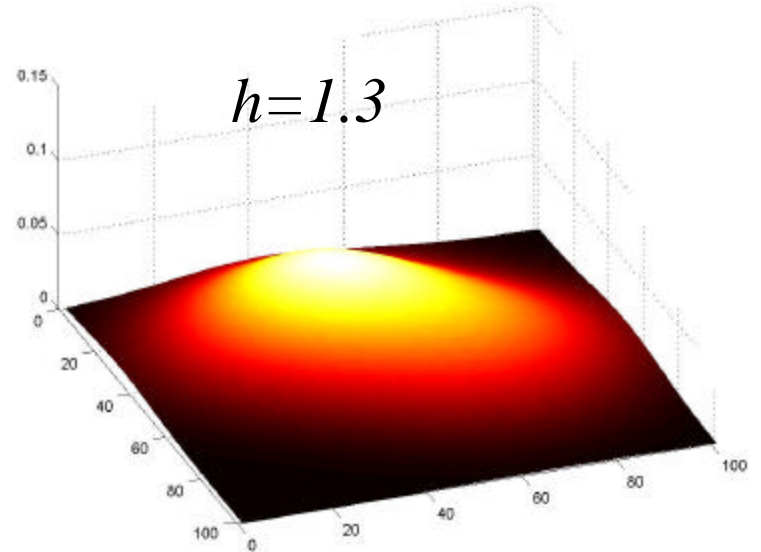
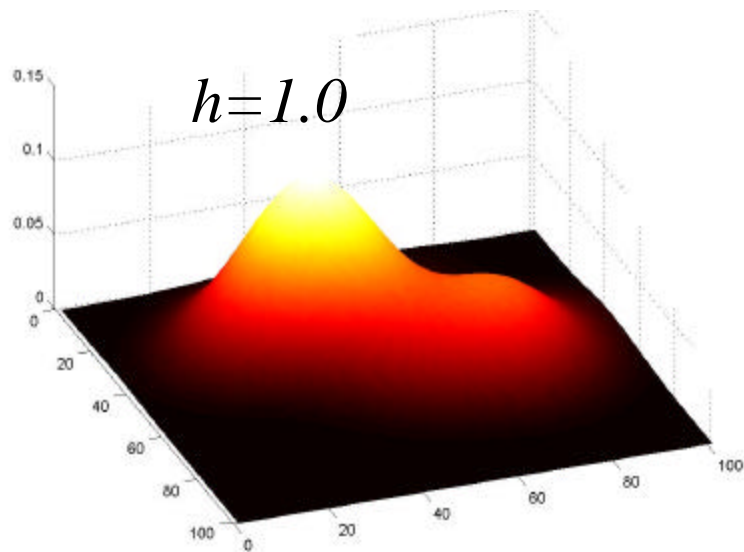
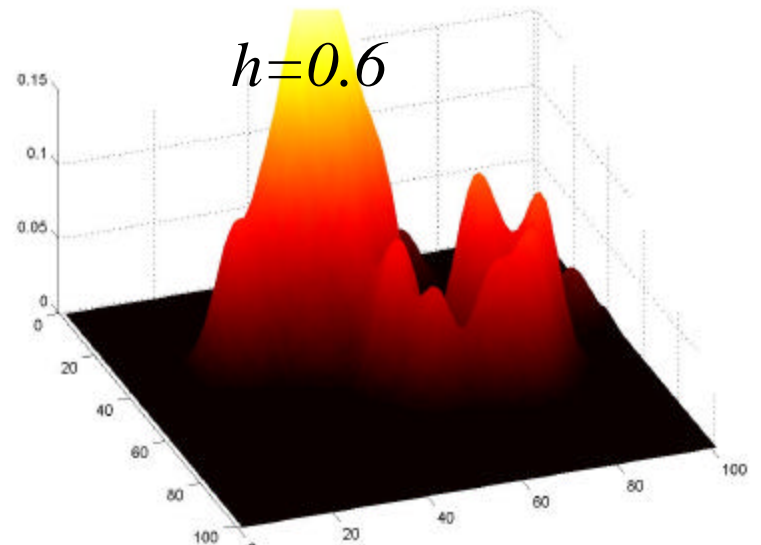
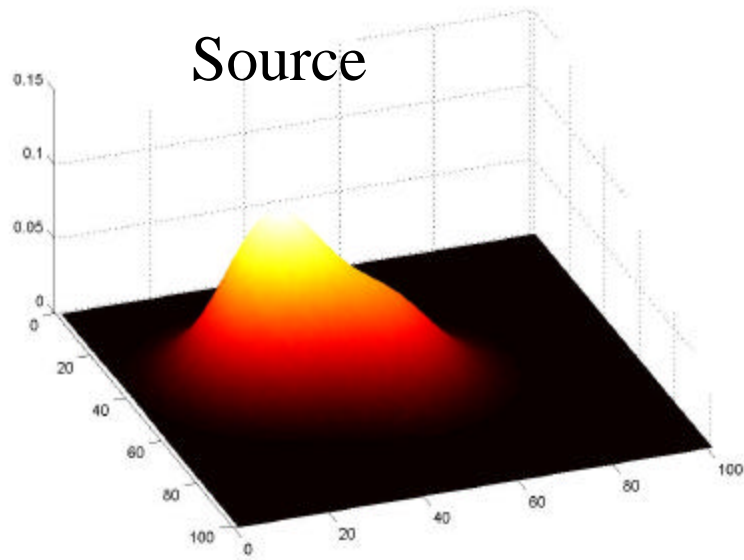
Eg: <loud cheer> a (spherical) Gaussian:

$$K(x, x^n) = \frac{1}{(\sqrt{2\mathbf{p}h})^d} \exp\left(-\frac{\|x - x^n\|^2}{2h^2}\right)$$

... so:

$$\tilde{p}(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(\sqrt{2\mathbf{p}h})^d} \exp\left(-\frac{\|x - x^n\|^2}{2h^2}\right)$$

Example



Some Insight

Interesting to look at expectation of the estimate with respect to all possible datasets:

$$E[\tilde{p}(x)] = E\left[\frac{1}{N} \sum_{n=1}^N K(x, x^n)\right] = E[K(x, x')]$$

$$= \int K(x - x') p(x') dx' \quad - \textit{convolution with true density}$$

That is:

$$\tilde{p}(x) = p(x) \quad \text{if} \quad K(x, x') = \mathbf{d}(x, x')$$

But not for the finite data set!

Conditions for Convergence

How small can we make h for a given N ?

$$\lim_{N \rightarrow \infty} h_N^d = 0$$

- It should go to 0

$$\lim_{N \rightarrow \infty} N h_N^d = \infty$$

- But slower than $1/N$

Based on the similar analysis of
variance of estimates

Eg:
$$h_N^d = h_1^d / \sqrt{N}$$

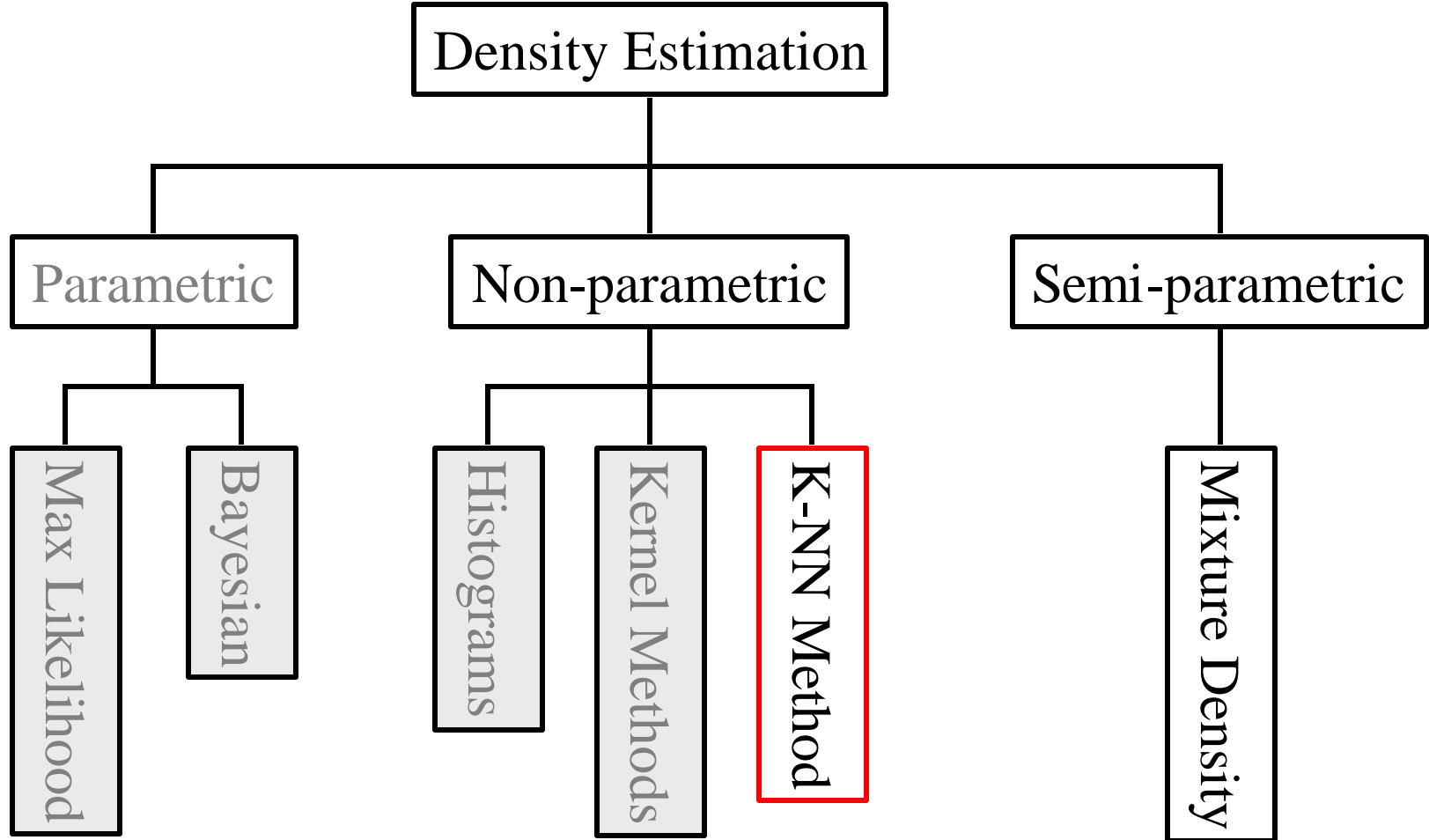
$$h_N^d = h_1^d / \log(N)$$

Note that the choice of h_1^d is still up to us.

Problems With Kernel Estimation

- Need to choose the width parameter, h
 - Can be chosen empirically
 - Can be adaptive, eg. $h_j = h d_{jk}$ – where d_{jk} the distance from x_j to k -th nearest neighbor
- Need to store all data to represent the density
 - Leads to Mixture Density Estimation

You Are Here



K-Nearest Neighbors

Recall that:

$$\tilde{p}(x) = \frac{k}{NV}$$

Now we fix k (typically $k = \sqrt{N}$) and expand V to contain k points

This is not a true density!

Eg.: choose $N=1$, $k=1$. Then:

$$\tilde{p}(x) = \frac{1}{1 \cdot \|x - x_1\|} \leftarrow \text{Oops!}$$

BUT it is useful for a number of theoretical and practical reasons.

K-NN Classification Rule

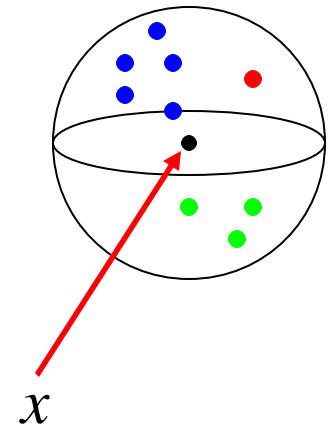
Let's try classification with K-NN density estimate

Data: N - total points
 N_j - points in class \mathbf{w}_j

Need to find the class label for a query, x

Expand a sphere from x to include K points

K - number of neighbors of x
 K_j - points of class \mathbf{w}_j among K



KNN Classification

Then class priors are given by: $p(\mathbf{w}_j) = \frac{N_j}{N}$

We can estimate conditional and marginal densities around any x :

$$p(x | \mathbf{w}_j) = \frac{K_j}{N_j V} \quad p(x) = \frac{K}{NV}$$

By Bayes rule: $p(\mathbf{w}_j | x) = \frac{K_j}{N_j V} \frac{N_j}{N} \frac{NV}{K} = \frac{K_j}{K}$

Then for *minimum error rate* classification:

$$C = \arg \max_j K_j$$

KNN Classification

Important theoretical result:

In the extreme case, $K=1$, it can be shown that:

N-sample error rate
for $P = \lim_{N \rightarrow \infty} P_N(\text{error})$

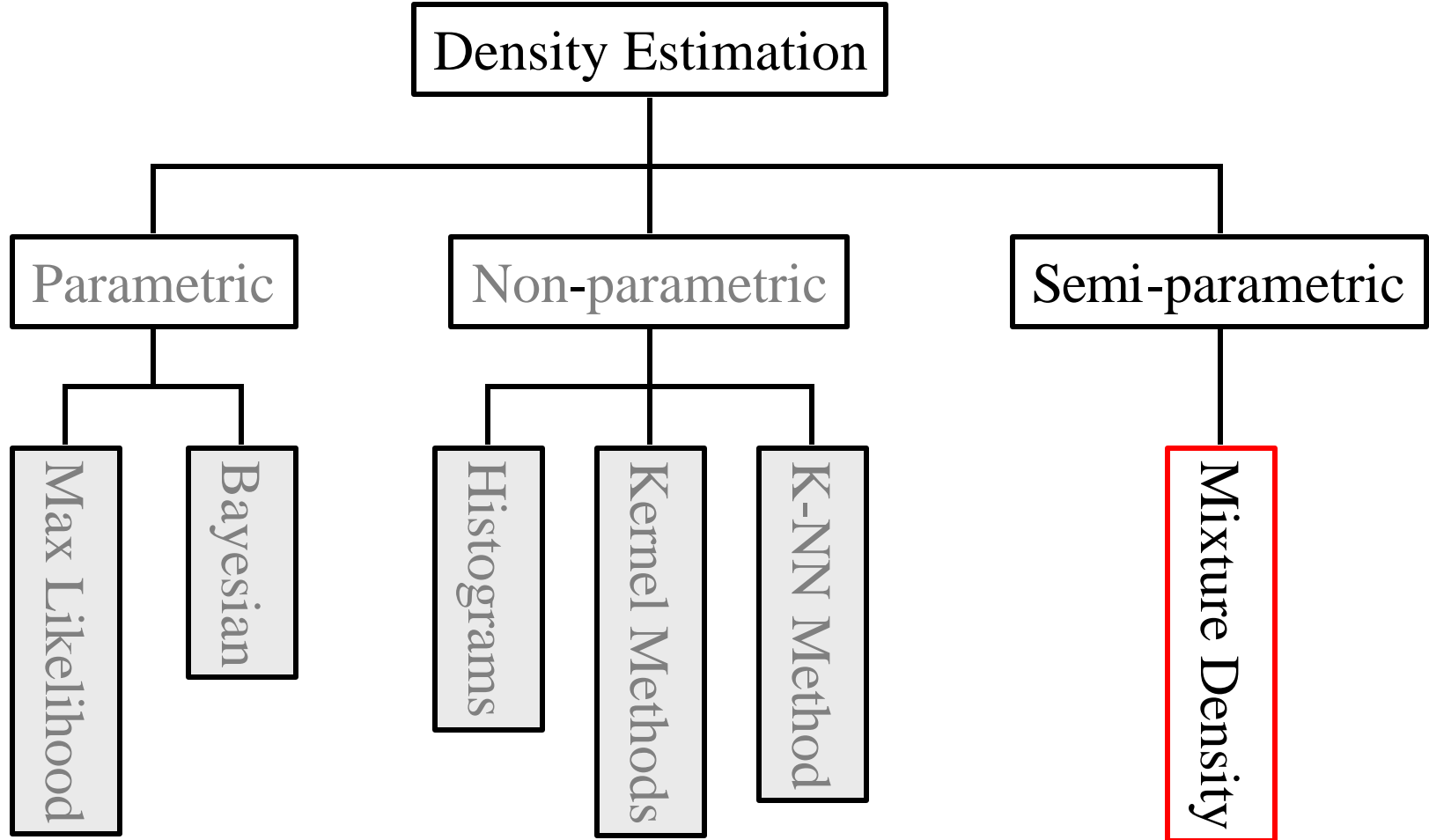
$$P^* \leq P \leq P^* \left(2 - \frac{c}{c-1} P^* \right)$$

That is, using just a single neighbor rule, the error rate is at most twice the Bayes error!!!

Problems with Non-parametric Methods

- Memory: need to store all data points
- Computation: need to compute distances to all data points every time
- Parameter choice: need to choose the smoothing parameter

You Are Here



Mixture Density Model

Mixture model – a linear combination of parametric densities

Number of components

$$p(x) = \sum_{j=1}^M p(x | j) P(j)$$

*Component
density*

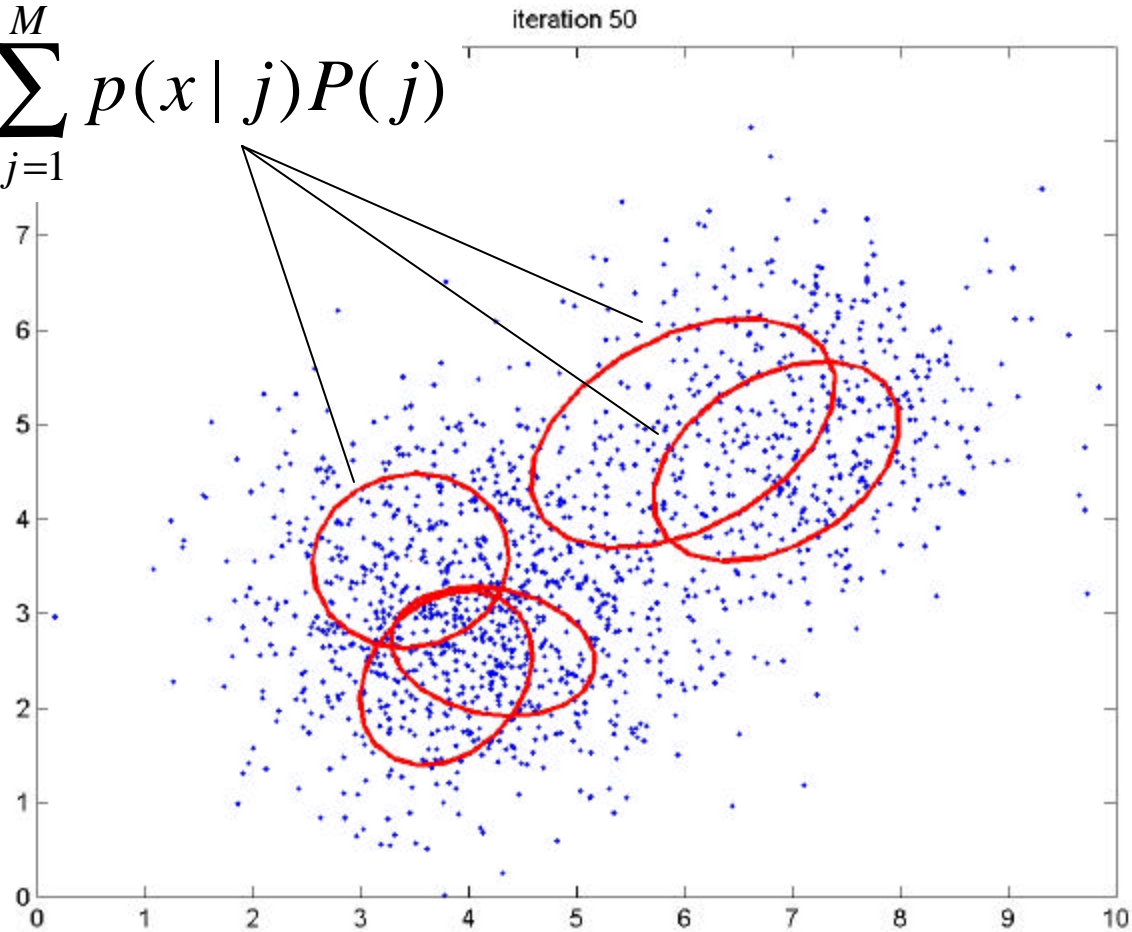
*Component
“prior”*

$$P(j) \geq 0, \quad \forall j \quad \text{and} \quad \sum_{j=1}^M P(j) = 1$$

Uses MUCH less “kernels” than kernel methods
Kernels are parametric densities, subject to estimation

Example

$$p(x) = \sum_{j=1}^M p(x | j) P(j)$$



Mixture Density

Using ML principle, the objective function is the *log-likelihood*:

$$l(\mathbf{q}) \equiv \log \prod_{n=1}^N p(x^n) = \sum_{n=1}^N \log \left\{ \sum_{j=1}^M p(x^n | j) P(j) \right\}$$

Differentiate w.r.t. parameters:

$$\begin{aligned} \nabla_{\mathbf{q}_j} l(\mathbf{q}) &= \sum_{n=1}^N \frac{\partial}{\partial \mathbf{q}_j} \log \left\{ \sum_{k=1}^M p(x^n | k) P(k) \right\} \\ &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^M p(x^n | k) P(k)} \frac{\partial}{\partial \mathbf{q}_j} p(x^n | j) P(j) \end{aligned}$$

Mixture Density

Again let's assume that $p(x/\omega)$ is a Gaussian

We need to estimate M priors, and M sets of means and covariances

$$\frac{\partial l(\mathbf{q})}{\partial \mathbf{m}_j} = \sum_{n=1}^N P(j | x^n) \left[\Sigma_j^{-1} (x^n - \hat{\mathbf{m}}_j) \right]$$

Setting it to 0 and solving for μ_j :

$$\hat{\mathbf{m}}_j = \frac{\sum_{n=1}^N P(j | x^n) x^n}{\sum_{n=1}^N P(j | x^n)}$$

- convex sum of all data

Mixture Density

Similarly for the covariances:

$$\frac{\partial l(\mathbf{q})}{\partial \mathbf{S}_j^2} = \sum_{n=1}^N P(j | x^n) \left[\hat{\mathbf{S}}_j^{-1} - \hat{\mathbf{S}}_j^{-1} (x^n - \hat{\mathbf{m}}_j)(x^n - \hat{\mathbf{m}}_j)^T \hat{\mathbf{S}}_j^{-1} \right]$$

Setting it to 0 and solving for Σ_i :

$$\hat{\mathbf{S}}_j = \frac{\sum_{n=1}^N P(j | x^n) (x^n - \hat{\mathbf{m}}_j)(x^n - \hat{\mathbf{m}}_j)^T}{\sum_{n=1}^N P(j | x^n)}$$

Mixture Density

A little harder for $P(j)$ – optimization is subject to constraints:

$$\sum_{j=1}^M P(j) = 1 \quad \text{and} \quad P(j) \geq 0, \forall j$$

Here is a trick to enforce the constraints:

$$P(j) = \frac{\exp(\mathbf{g}_j)}{\sum_{k=1}^M \exp(\mathbf{g}_k)}$$

$$\frac{\partial P(i)}{\partial \mathbf{g}_j} = \mathbf{d}(i-j)P(j) - P(i)P(j)$$

Mixture Density

Using the chain rule:

$$\begin{aligned}\nabla_{\mathbf{g}_j} l(\mathbf{q}) &= \sum_{k=1}^M \frac{\partial l(\mathbf{q})}{\partial P(k)} \frac{\partial P(k)}{\partial \mathbf{g}_j} = \sum_{k=1}^M \sum_{n=1}^N \frac{p(x^n | k)}{P(x)} (\mathbf{d}_{jk} P(j) - P(j)P(k)) \\ &= \sum_{n=1}^N \left\{ \frac{p(x^n | j)}{P(x)} P(j) - \sum_{k=1}^M \frac{p(x^n | k)}{P(x)} P(j)P(k) \right\} \\ &= \sum_{n=1}^N \left\{ P(j | x^n) - P(j) \sum_{k=1}^M p(k | x^n) \right\} = \sum_{n=1}^N \{ P(j | x^n) - P(j) \} = 0\end{aligned}$$

The last expression gives the value at the extremum:

$$P(j) = \frac{1}{N} \sum_{n=1}^N P(j | x^n)$$

Mixture Density

What's the problem?

$$P(j) = \frac{1}{N} \sum_{n=1}^N P(j | x^n)$$

$$\hat{\mathbf{m}}_j = \frac{\sum_{n=1}^N P(j | x^n) x^n}{\sum_{n=1}^N P(j | x^n)}$$

$$\hat{\mathbf{S}}_j = \frac{\sum_{n=1}^N P(j | x^n) (x^n - \hat{\mathbf{m}}_j)(x^n - \hat{\mathbf{m}}_j)^T}{\sum_{n=1}^N P(j | x^n)}$$

We can't compute these directly!

Solution – EM algorithm. We will study it in Clustering.

You Are Here

