The EP Algorithm

- Given: a joint distribution over data and variables $p(\mathcal{D}, \theta) = \prod_{i=1}^{M} f_i(\theta)$
- Goal: approximate the posterior $p(\theta \mid D)$ with q
- Initialize all approximating factors $\tilde{f}_i(\theta)$
- Initialize the posterior approximation $q(\theta) \propto \prod \tilde{f}_i(\theta)$
- Do until convergence:
 - choose a factor $\tilde{f}_j(\boldsymbol{\theta})$
 - remove the factor from q by division: $q^{\setminus j}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_i(\boldsymbol{\theta})}$





The EP Algorithm

• find $q^{\rm new}$ that minimizes

$$\operatorname{KL}\left(\frac{f_j(\theta)q^{\setminus j}(\boldsymbol{\theta})}{Z_j}\Big|q^{\operatorname{new}}(\boldsymbol{\theta})\right)$$

using moment matching, including the zero-th moment: $\int_{C} dx$

$$Z_j = \int q^{\setminus j}(\boldsymbol{\theta}) f_j(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

evaluate the new factor

$$\widetilde{f}_j(\boldsymbol{\theta}) = Z_j \frac{q^{\text{new}}(\boldsymbol{\theta})}{q^{\setminus j}(\boldsymbol{\theta})}$$

• After convergence, we have $p(\mathcal{D}) \approx \int \prod_{j} \tilde{f}_{j}(\theta) d\theta$



Example



yellow: original distribution

red: Laplace approximation

green: global variation

blue: expectation-propagation



The Clutter Problem



Aim: fit a multivariate Gaussian into data in the presence of background clutter (also Gaussian) p(x | θ) = (1 - w)N(x | θ, I) + wN(x | 0, aI)
The prior is Gaussian: p(θ) = N(θ | 0, bI)



The Clutter Problem

The joint distribution for $\mathcal{D}_{N} = (\mathbf{x}_{1}, \dots, \mathbf{x}_{N})$ is $p(\mathcal{D}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{n=1}^{N} p(\mathbf{x}_{n} \mid \boldsymbol{\theta})$

this is a mixture of 2^N Gaussians! This is intractable for large *N*. Instead, we approximate it using a spherical Gaussian:

$$q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{m}, vI) = \tilde{f}_0(\boldsymbol{\theta}) \prod_{n=1}^N \tilde{f}_n(\boldsymbol{\theta})$$

the factors are (unnormalized) Gaussians:

$$\tilde{f}_0(\boldsymbol{\theta}) = p(\boldsymbol{\theta}) \qquad \tilde{f}_n(\boldsymbol{\theta}) = s_n \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{m}_n, v_n I)$$



EP for the Clutter Problem

- First, we initialize $\tilde{f}_n(\theta) = 1$, i.e. $q(\theta) = p(\theta)$
- Iterate:
 - Remove the current estimate of $f_n(\theta)$ from q by division of Gaussians:

$$q_{-n}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_n(\boldsymbol{\theta})}$$



EP for the Clutter Problem

- First, we initialize $\tilde{f}_n(\theta) = 1$, i.e. $q(\theta) = p(\theta)$
- Iterate:
 - Remove the current estimate of $\tilde{f}_n(\theta)$ from q by division of Gaussians:

$$q_{-n}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_n(\boldsymbol{\theta})} \qquad q_{-n}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{m}_{-n}, v_{-n}I)$$

Compute the normalization constant:

$$Z_n = \int q_{-n}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

= $(1 - w) \mathcal{N}(\mathbf{x}_n \mid \mathbf{m}_{-n}, (v_{-n} + 1)I) + w \mathcal{N}(\mathbf{x}_n \mid \mathbf{0}, aI)$



EP for the Clutter Problem

- First, we initialize $\tilde{f}_n(\theta) = 1$, i.e. $q(\theta) = p(\theta)$
- Iterate:
 - Remove the current estimate of $\tilde{f}_n(\theta)$ from q by division of Gaussians:

$$q_{-n}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_n(\boldsymbol{\theta})} \qquad q_{-n}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{m}_{-n}, v_{-n}I)$$

Compute the normalization constant:

$$Z_n = \int q_{-n}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

- Compute mean and variance of $q^{\text{new}} = q_{-n}(\theta) f_n(\theta)$
- Update the factor $\tilde{f}_n(\theta) = Z_n \frac{q^{\text{new}}(\theta)}{q_{-n}(\theta)}$



A 1D Example



- blue: true factor $f_n(\theta)$
- red: approximate factor $\tilde{f}_n(\theta)$
- green: cavity distribution $q_{-n}(\theta)$

The form of $q_{-n}(\theta)$ controls the range over which $\tilde{f}_n(\theta)$ will be a good approximation of $f_n(\theta)$



Summary

- Variational Inference uses approximation of functions so that the KL-divergence is minimal
- In mean-field theory, factors are optimized sequentially by taking the expectation over all other variables
- Variational inference for GMMs reduces the risk of overfitting; it is essentially an EM-like algorithm
- Expectation propagation minimizes the reverse KL-divergence of a single factor by moment matching; factors are in the exp. family





Computer Vision Group Prof. Daniel Cremers Technische Universität München

10. Sampling Methods

Sampling Methods

Sampling Methods are widely used in Computer Science

- as an approximation of a deterministic algorithm
- to represent uncertainty without a parametric model
- to obtain higher computational efficiency with a small approximation error
- Sampling Methods are also often called Monte Carlo Methods
- Example: Monte-Carlo Integration
 - Sample in the bounding box
 - Compute fraction of inliers
 - Multiply fraction with box size





Non-Parametric Representation

Probability distributions (e.g. a robot's belief) can be represeted:

- Parametrically: e.g. using mean and covariance of a Gaussian
- Non-parametrically: using a set of hypotheses (samples) drawn from the distribution

Advantage of non-parametric representation:

 No restriction on the *type* of distribution (e.g. can be multi-modal, non- Gaussian, etc.)



Non-Parametric Representation



The more samples are in an interval, the higher the probability of that interval

But:

How to draw samples from a function/distribution?



Sampling from a Distribution

There are several approaches:

- Probability transformation
 - Uses inverse of the c.d.f h
- Rejection Sampling
- Importance Sampling
- MCMC

But:

Probability transformation:

- Sample uniformly in [0,1]/
- Transform using h⁻¹



Requires calculation of h and its inverse

Machine Learning for Computer Vision



Rejection Sampling

1. Simplification:

- Assume p(z) < 1 for all z
- Sample z uniformly
- Sample c from [0,1]

If f(z) > c :
 keep the sample
 otherwise:

reject the sample





Rejection Sampling

2. General case:

Assume we can evaluate $p(z) = \frac{1}{Z_n} \tilde{p}(z)$ (unnormalized)

- Find proposal distribution q
 - Easy to sample from q
- Find k with $kq(z) \ge \tilde{p}(z)$
- Sample from q
- Sample uniformly from [0,kq(z₀)]
- Reject if $u_0 > \tilde{p}(z_0)$



But: Rejection sampling is inefficient.



Importance Sampling

- Idea: assign an importance weight w to each sample
- With the importance weights, we can account for the "differences between p and q "

w(x) = p(x)/q(x)

- p is called target
- q is called proposal (as before)





Importance Sampling

- Explanation: The prob. of falling in an interval A is the area under p
- This is equal to the expectation of the indicator function $I(x \in A)$

$$E_p[I(z \in A)] = \int p(z)I(z \in A)dz$$

$$\begin{array}{c} p(z) \\ A \end{array}$$



Importance Sampling

- Explanation: The prob. of falling in an interval A is the area under p
- This is equal to the expectation of the indicator function $I(x \in A)$

$$E_p[I(z \in A)] = \int p(z)I(z \in A)dz$$

$$\begin{array}{c|c} p(z) \\ A \end{array} \end{array}$$

 \sim ()

 $= \int \frac{p(z)}{q(z)} q(z) I(z \in A) dz = E_q[w(z)I(z \in A)]$ Requirement: $p(x) > 0 \Rightarrow q(x) > 0$

Approximation with samples drawn from q: $E_q[w(z)I(z \in A)] \approx \frac{1}{L} \sum_{l=1}^{L} w(z_l)I(z_l \in A)$



The Particle Filter

- Non-parametric implementation of Bayes filter
- Represents the belief (posterior) $Bel(x_t)$ by a set of random state samples.
- This representation is approximate.
- Can represent distributions that are **not Gaussian**.
- Can model non-linear transformations.

Basic principle:

- Set of state hypotheses ("particles")
- Survival-of-the-fittest



The Bayes Filter Algorithm (Rep.)

$$Bel(x_t) = \eta \ p(z_t \mid x_t) \ \int p(x_t \mid u_t, x_{t-1}) Bel(x_{t-1}) dx_{t-1}$$

Algorithm Bayes_filter (Bel(x), d)

1. if d is a sensor measurement z then

$$2. \quad \eta = 0$$

3. for all x do

4.
$$\operatorname{Bel}'(x) \leftarrow p(z \mid x) \operatorname{Bel}(x)$$

5.
$$\eta \leftarrow \eta + \operatorname{Bel}'(x)$$

- 6. for all x do $\operatorname{Bel}'(x) \leftarrow \eta^{-1} \operatorname{Bel}'(x)$
- 7. else if d is an action u then
- 8. for all x do $Bel'(x) \leftarrow \int p(x \mid u, x')Bel(x')dx'$
- 9. return $\operatorname{Bel}'(x)$



Mathematical Description

Set of weighted samples:



The samples represent the probability distribution:

$$p(x) = \sum_{i=1}^{M} w_t^{[i]} \cdot \delta_{x_t^{[i]}}(x)$$
 Point mass
distribution
("Dirac")



The Particle Filter Algorithm





Localization with Particle Filters

- Each particle is a potential pose of the robot
- Proposal distribution is the motion model of the robot (prediction step)
- The observation model is used to compute the importance weight (correction step)
- Randomized algorithms are usually called Monte Carlo algorithms, therefore we call this:

Monte-Carlo Localization



A Simple Example



- The initial belief is a uniform distribution (global localization).
- This is represented by an (approximately) uniform sampling of initial particles.



Sensor Information



The sensor model $p(z_t | x_t^{[m]})$ is used to compute the new importance weights:

$$w_t^{[m]} \leftarrow p(z_t \mid x_t^{[m]})$$



Robot Motion



After resampling and applying the motion model $p(x_t \mid u_t, x_{t-1}^{[m]})$ the particles are distributed more densely at three locations.





Sensor Information



Again, we set the new importance weights equal to the sensor model.

$$w_t^{[m]} \leftarrow p(z_t \mid x_t^{[m]})$$



Robot Motion



Resampling and application of the motion model: One location of dense particles is left. **The robot is localized.**



A Closer Look at the Algorithm...





Sampling from Proposal

This can be done in the following ways:

- Adding the motion vector to each particle directly (this assumes perfect motion) $[m]_{t,x_t=1}$
- Sampling from the motion model , e.g. for a 2D motion with translation velocity v and rotation velocity w we have: $p(x_t | u_t, x_{t-1}^{[m]})$

$$\mathbf{u}_{t} = \begin{pmatrix} v_{t} \\ w_{t} \end{pmatrix} \qquad \mathbf{x}_{t} = \begin{pmatrix} x_{t} \\ y_{t} \\ \theta_{t} \end{pmatrix} \qquad \text{Orientation}$$



Motion Model Sampling (Example)





Computation of Importance Weights

Computation of the sample weights:

- Proposal distribution: $g(x_t^{[m]}) = p(x_t^{[m]} | u_t, x_{t-1}^{[m]}) \text{Bel}(x_{t-1}^{[m]})$ (we sample from that using the motion model)
- Target distribution (new belief): $f(x_t^{[m]}) = \text{Bel}(x_t^{[m]})$ (we can not directly sample from that \rightarrow importance sampling)
- Computation of importance weights:

$$\boxed{w_t^{[m]}} = \frac{f(x_t^{[m]})}{g(x_t^{[m]})} \propto \frac{p(z_t \mid x_t^{[m]})p(x_t^{[m]} \mid u_t, x_{t-1}^{[m]})\operatorname{Bel}(x_{t-1}^{[m]})}{p(x_t^{[m]} \mid u_t, x_{t-1}^{[m]})\operatorname{Bel}(x_{t-1}^{[m]})} = \boxed{p(z_t \mid x_t^{[m]})}$$



Proximity Sensor Models

• How can we obtain the sensor model $p(z_t \mid x_t^{[m]})$? Sensor Calibration:



Laser sensor

Dr. Rudolph Triebel

Computer Vision Group



Resampling

- Given: Set $\bar{\mathcal{X}}_t$ of weighted samples.
- Wanted : Random sample, where the probability of drawing x_i is equal to w_i.
- Typically done M times with replacement to generate new safe $m = \operatorname{Pt} M$ to do draw *i* with prob. $\propto w_t^{[i]}$ $\chi_t \leftarrow \chi_t \cup x_t^{[i]}$ χ_t



Resampling





Standard n-times sampling results in high variance
This requires more particles
O(nlog n) complexity

- Instead: low variance sampling only samples once
- Linear time complexity
- Easy to implement



Sample-based Localization (sonar)





Initial Distribution





After Ten Ultrasound Scans





After 65 Ultrasound Scans





Estimated Path





Kidnapped Robot Problem

The approach described so far is able to

- track the pose of a mobile robot and to
- globally localize the robot.
- How can we deal with localization errors (i.e., the kidnapped robot problem)?
- Idea: Introduce uniform samples at every resampling step
- This adds new hypotheses and reduces the





Summary

- There are mainly 4 different types of sampling methods: Transformation method, rejections sampling, importance sampling and MCMC
- Transformation only rarely applicable
- Rejection sampling is often very inefficient
- Importance sampling is used in the particle filter which can be used for robot localization
- An efficient implementation of the resampling step is the low variance sampling



