

# 1 System Modelling

Continuous time (linear):

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \mathbf{x}(0) = \mathbf{x}_0 \quad \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$$

Continuous time (nonlinear):

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)), \mathbf{x}(0) = \mathbf{x}_0 \quad \mathbf{y}(t) = \mathbf{h}(\mathbf{x}(t), \mathbf{u}(t))$$

Discrete time (linear):

$$\mathbf{x}_{k+1} = \mathbf{A}_d\mathbf{x}_k + \mathbf{B}_d\mathbf{u}_k \quad \mathbf{y}_k = \mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k$$

With a zero-order hold:

$$\begin{bmatrix} \mathbf{A}_d & \mathbf{B}_d \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \exp \left( \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} h \right) \quad \mathbf{A}_d = e^{\mathbf{A}h}, \mathbf{B}_d = \int_0^h e^{\mathbf{A}\tau'} d\tau' \mathbf{B}$$

## 2 Numerical Stability

Well conditioned problem:  $\text{Cond}_x$  small or  $\text{cond}_x \leq 1$

$$\text{Cond}_x = \left| \frac{\Delta \tilde{y}}{\Delta x} \right| \approx \left| \frac{df}{dx} \right| \quad \text{cond}_x = \left| \frac{\delta \tilde{y}}{\delta x} \right| \approx |K_x| = \left| \frac{df}{dx} \cdot \frac{x}{f(x)} \right|$$

$\varphi$  is order  $p$  or  $O(\Delta^p)$  if:  $\tilde{y} = \varphi(x, \Delta) - f(x) \propto \Delta^p$

Consistency:  $\lim_{\Delta \rightarrow 0} \varphi(x, \Delta) = f(x)$       Stability:  $\left| \frac{\Delta \tilde{y}}{\Delta x} \right| \approx \left| \frac{d\varphi}{dx} \right| < 1$

$\hat{y} = \varphi(\hat{x}, \Delta) = \varphi(\hat{x} + \Delta x)$  converges to  $f(x)$  for  $\Delta \rightarrow 0$  if it's consistent and at least marginally stable.

$\hat{\mathbf{x}} = \mathbf{A}\hat{\mathbf{x}}$  well-conditioned/stable when  $\text{Re}(\lambda) < 0$  for all  $\lambda$ .

$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k$  well-conditioned/stable when  $|\lambda| < 1$  for all  $\lambda$ .

## 3 Numerical Root Finding

$f$  is Lipschitz continuous if

$$\exists c \text{ s.t. } \forall x, z \in \mathbb{R}, |f(x) - f(z)| \leq c|x - z|$$

An algorithm has convergence rate/order  $r$  if:  $\lim_{k \rightarrow \infty} \frac{|E_{k+1}|}{|E_k|^r} = C$

**Bisection:** order 1; requires continuity; have  $l_k, r_k$  such that  $\text{sgn } f(l_k) \neq \text{sgn } f(r_k)$ , check  $c = \frac{l_k + r_k}{2}$  each iteration.

**Fixed-point Iteration:** order 2 if  $g'(x^*) \approx 0, 1$  otherwise; requires update function  $g(x^*) = x^*$  for Lipschitz  $g$  with  $c < 1$  for  $x \approx x^*$ . To find  $g(x)$ , rearrange  $f(x^*) = 0$  into the form  $x^* = g(x^*)$ . Start with guess, update as  $x_{k+1} = g(x_k)$ .

**Newton's Method:** order 2; requires continuous differentiability; equivalent to fixed-point iteration with  $g(x) = x - \frac{f(x)}{f'(x)}$ . Start with guess, update as

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}. \text{ Exact for linear } f \text{ and suffers for highly nonlinear } f.$$

**Secant Method:** order 1.6; Newton with  $f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$ .

## 4 Numerical Integration and Differentiation

**Midpoint/Trapezoidal Rule:** order 2, approx.  $f$  as const./linear

$$\int_{x_i}^{x_{i+1}} f(x) dx \approx (x_{i+1} - x_i) \cdot f\left(\frac{x_{i+1} + x_i}{2}\right) \approx (x_{i+1} - x_i) \cdot \frac{f(x_{i+1}) + f(x_i)}{2}$$

**Simpson's Rule:** order 4, approx.  $f$  as quadratic

$$(x_{i+1} - x_i) \cdot \frac{f(x_{i+1}) + 4f\left(\frac{x_{i+1} + x_i}{2}\right) + f(x_i)}{6}$$

**Differentiation:** Forward (order 1), Backward (1), Centered (2) Difference

$$f'(x) \approx \frac{f(x + \Delta) - f(x)}{\Delta} \approx \frac{f(x) - f(x - \Delta)}{\Delta} \approx \frac{f(x + \Delta) - f(x - \Delta)}{2\Delta}$$

## 5 Numerical ODE Solving

**Forward Euler's Method:** order 1 (global); explicit; conditionally stable if  $|1 + h\lambda| < 1$ :  $x_{k+1} = x_k + hf_k$

**Backward Euler's Method:** order 1 (global); implicit; conditionally stable if  $|1 - h\lambda| > 1$ :  $x_{k+1} = x_k + hf_{k+1}$

**Trapezoidal Method:** order 2 (global); implicit; conditionally stable if  $\left| \frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} \right| < 1$ :  $x_{k+1} = x_k + \frac{h}{2}(f_{k+1} + f_k)$

**Heun's Method (RK2):** order 2 (global); explicit; conditionally stable if  $-4 < 2h\lambda + h^2\lambda^2 < 0$ :  $x_{k+1} = x_k + \frac{h}{2}(f(x_k + hf_k) + f_k)$

**Fourth-Order Runge-Kutta (RK4):** order 4 (global); explicit.

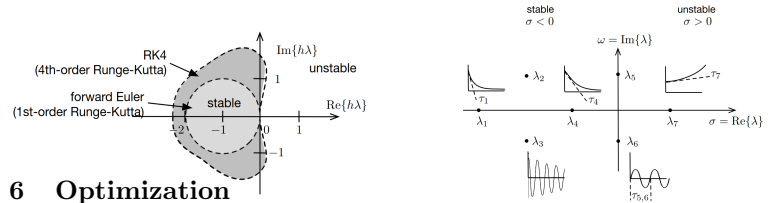
$$x_{k+1} = x_k + \frac{h}{6}(k_1 + 2k_2 + 2k_3 + k_4)$$

$$k_1 = f(x_k), \quad k_2 = f\left(x_k + \frac{1}{2}hk_1\right) \quad k_3 = f\left(x_k + \frac{1}{2}hk_2\right) \quad k_4 = f(x_k + hk_3)$$

*Time constant:* for  $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$  with complex eigenvalues  $\lambda_i$ :

$$\tau_i = \left\{ \frac{1}{|\text{Re}(\lambda_i)|}, \frac{2\pi}{|\text{Im}(\lambda_i)|} \right\} \quad \gamma = \frac{\tau_{\max}}{\tau_{\min}} \text{ stiff if } > 10^3$$

Rule of thumb: simulate  $T = 5\tau_{\max}$  if system is stable; use step size  $h = \min \left\{ \frac{\tau_{\min}}{10}, \frac{T}{200} \right\}$  or number of steps  $k = \max \{ 50\gamma, 200 \}$ ; plot with step size  $H = \frac{T}{200} = \frac{\tau_{\max}}{40}$ ; for stiff systems use a variable-step solver with initial step  $\frac{\tau_{\min}}{10}$ . Use adaptive step size solvers for stiff problems.



## 6 Optimization

Formulation: minimize  $f(\mathbf{x})$ , subject to  $\mathbf{g}(\mathbf{x}) = \mathbf{0}$  (equality constraints) and  $\mathbf{h}(\mathbf{x}) \geq \mathbf{0}$  (inequality constraints);  $f: \mathbb{R}^n \mapsto \mathbb{R}, \mathbf{g}: \mathbb{R}^n \mapsto \mathbb{R}^m, \mathbf{h}: \mathbb{R}^n \mapsto \mathbb{R}^p$ .  $\mathbf{x}^*$  is a *global minimum* if  $\forall \mathbf{x} \in \mathbb{R}^n, f(\mathbf{x}^*) \leq f(\mathbf{x})$ ; or *local minimum* if  $\exists \varepsilon > 0$  s.t.  $\forall \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x} - \mathbf{x}^*\|_2 < \varepsilon \implies f(\mathbf{x}^*) \leq f(\mathbf{x})$ . All minima of  $f$  are *stationary points*, where  $\vec{\nabla} f(\mathbf{x}) = \mathbf{0}$ ; to find the global minimum, check all stationary points and boundaries.

$$\mathbf{H}_f(\mathbf{x}) = \text{matrix}_{i,j} \left[ \frac{\partial^2 f}{\partial x_i \partial x_j} \right]$$

Local minimum if  $\mathbf{H}_f(\mathbf{x}^*) > 0$ , maximum if  $\mathbf{H}_f < 0$ , saddle if  $\mathbf{H}_f = 0$ .

A *feasible point* is any  $\mathbf{x}$  satisfying all constraints; the *feasible set* contains all feasible points. A *critical point* is a local maximum, minimum or saddle point in the feasible set.

*Lagrange Multipliers:* for equality constraints  $\mathbf{g}(\mathbf{x}) = \mathbf{0}$ :

$$\min_{\mathbf{x}, \boldsymbol{\lambda}} \Lambda(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$$

*Karush-Kuhn-Tucker (KKT) Conditions:*  $\mathbf{x}^*$  is a critical point when there exists  $\boldsymbol{\lambda} \in \mathbb{R}^m$  and  $\boldsymbol{\mu} \in \mathbb{R}^p$  such that:

- Stationarity:  $\vec{\nabla} f(\mathbf{x}^*) - \sum_i \lambda_i \vec{\nabla} g_i(\mathbf{x}^*) - \sum_j \mu_j \vec{\nabla} h_j(\mathbf{x}^*) = \mathbf{0}$
- Primal feasibility:  $\mathbf{g}(\mathbf{x}^*) = \mathbf{0}$  and  $\mathbf{h}(\mathbf{x}^*) \geq \mathbf{0}$
- Complementary slackness:  $\forall j, \mu_j h_j(\mathbf{x}^*) = 0$
- Dual feasibility:  $\forall j, \mu_j \geq 0$

$f$  is *convex* if  $\mathbf{H}_f > 0$  for all  $\mathbf{x}$ , or

$$\forall \mathbf{x}_1 \neq \mathbf{x}_2, \forall \alpha \in (0, 1), f((1 - \alpha)\mathbf{x}_1 + \alpha\mathbf{x}_2) \leq (1 - \alpha)f(\mathbf{x}_1) + \alpha f(\mathbf{x}_2)$$

A set  $S$  is *convex* if:  $\forall \mathbf{x}, \mathbf{y} \in S, \alpha \in [0, 1], \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in S$

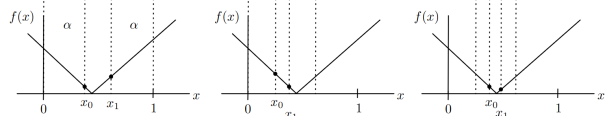
or  $S = \{ \mathbf{x} \mid f(\mathbf{x}) \leq c \}$  for some convex  $f(\mathbf{x})$ .

If the objective  $f$  and feasible set are both convex, the problem is convex and has a unique minimum.

## 7 Numerical Optimization Algorithms

$f$  is *unimodal* over  $[a, b]$  if there exists  $x^* \in [a, b]$  s.t.  $f$  is decreasing over  $[a, x^*]$  and increasing over  $[x^*, b]$ .

**Golden Section Search:** 1D, unconstrained; order 1; requires unimodality. Rescale to  $[0, 1]$ , choose  $x_0 = \alpha, x_1 = 1 - \alpha$  for  $0 < \alpha < 1/2$ ; evaluate  $f(x_0), f(x_1)$ , discard larger side, rescale and repeat. Using  $1 - \alpha = \frac{1}{2}(\sqrt{5} - 1)$  allows using  $x_0$  as the next  $x_1$ .



**Gradient Descent:** unconstrained; requires twice differentiable. Let  $g(\alpha) = f(\mathbf{x}_k - \alpha \vec{\nabla} f(\mathbf{x}_k)^T)$ , find  $\alpha^*$  minimizing  $g$ , update  $x_{k+1} = x_k - \alpha^* \vec{\nabla} f(\mathbf{x}_k)^T$ . Optimizing  $g$  is expensive so fixed step size can be used. Suffers from poor conditioning of  $f$ .

**Newton's Method:** unconstrained; requires twice differentiable. Iterate as  $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}_f^{-1}(\mathbf{x}_k) \vec{\nabla} f(\mathbf{x}_k)^T$ ; or in 1D,  $x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}$ . Gauss-Newton approximates  $\mathbf{H}_f$  using first derivatives (mix of gradient descent and Newton). Levenberg-Marquardt applies adaptive regularization for nearly-singular Hessians.

**Sequential Quadratic Programming (SQP):** constrained. Iteratively solves a series of simpler, less constrained approximations of the problem.  $f$  is replaced by a quadratic approximation and the constraints linearized. Similar to Newton's method; only converges if initial guess is good.

**Barrier Methods:** constrained. Constraints are turned into penalties on the objective. Define new objective as  $f'(x) = f(x) + \rho \frac{1}{h(x)}$  where weight  $\rho$  is increased to satisfy constraints, decreased for more accuracy.

## 8 Linear Algebra

$\mathbf{A}^T \mathbf{A}$  is positive definite if  $\mathbf{A}$  is full rank, semi-definite otherwise.

Orthogonal  $\mathbf{Q}$  rotates, preserves angles, and  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ .

$\mathbf{A}\mathbf{x} = \mathbf{b}$  can be solved by factoring  $\mathbf{A} = \mathbf{L}\mathbf{U}$  and solving  $\mathbf{L}\mathbf{y} = \mathbf{b}, \mathbf{U}\mathbf{x} = \mathbf{y}$ .

A general *vector norm* satisfies:  $\|\mathbf{x}\| = 0 \iff \mathbf{x} = \mathbf{0}$ ;  $\|c\mathbf{x}\| = |c|\|\mathbf{x}\|$ ;  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ . The  $p$ -norm for  $p \geq 1$  is convex:

$$\|\mathbf{x}\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}}$$

2-norm is Euclidean;  $\infty$ -norm is  $\|\mathbf{x}\|_\infty = \max\{|x_1|, |x_2|, \dots, |x_n|\}$ .

The *matrix norm* induced by a vector norm is

$$\|\mathbf{A}\| = \max \{ \|\mathbf{Ax}\| \mid \|\mathbf{x}\| = 1 \} = \max_{\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|}$$

1-norm: max col. sum; 2-norm/spectral radius: "largest eigenvalue" of  $\mathbf{A}$

$$\|\mathbf{A}\|_1 = \max_{i \leq j \leq n} \sum_{i=1}^m |a_{ij}| \quad \|\mathbf{A}\|_2 = \max \left\{ \sqrt{\lambda} \mid \exists \mathbf{x} \in \mathbb{R}^n \text{ s.t. } \mathbf{A}^T \mathbf{Ax} = \lambda \mathbf{x} \right\}$$

Frobenius norm: always  $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F$ ;  $\infty$ -norm: maximum row sum

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{tr } \mathbf{A}^T \mathbf{A}} \quad \|\mathbf{A}\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$$

The *condition number* of  $\mathbf{A}$  with respect to a norm is  $\text{cond } \mathbf{A} = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  or  $\infty$  if  $\mathbf{A}$  is non-invertible. For solving  $\mathbf{Ax} = \mathbf{b}$ , this describes how error in  $\mathbf{A}$  and  $\mathbf{b}$  propagates to  $\mathbf{x}$ .

**Eigendecomposition:**  $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1}$  where  $\mathbf{Av} = \lambda\mathbf{v}$ .

**Singular Value Decomposition:**  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  where  $\mathbf{V} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{U} \in \mathbb{R}^{m \times m}$  are orthogonal,  $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$  is rectangular diagonal with singular values. Singular values and vectors are related by  $\mathbf{Av}_i = \sigma_i \mathbf{u}_i$ .

$\sigma_i = \sqrt{\lambda_i(\mathbf{A}^T \mathbf{A})} = \sqrt{\lambda_i(\mathbf{A}\mathbf{A}^T)}$ , total  $\max\{m, n\}$  singular values, nonnegative and real, in decreasing order with  $k = \text{rank } \mathbf{A}$  nonzero singular values.

$\mathbf{V}$  has  $k$  eigenvectors of  $\mathbf{A}^T \mathbf{A}$  and  $n - k$  orthonormal vectors from  $\ker(\mathbf{A}^T \mathbf{A})$ ;  $\mathbf{U}$  has  $k$  eigenvectors of  $\mathbf{A}\mathbf{A}^T$  and  $m - k$  orthonormal vectors from  $\ker(\mathbf{A}\mathbf{A}^T)$ .

The rank  $j \leq k$  approximation of  $\mathbf{A}$  is  $\sum_{i=1}^j \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ .

## 9 Probability Definitions and Theorems

The *cumulative distribution function* (CDF) of  $f_x(x)$  is

$$F_x(x) = \sum_{\tilde{x} = -\infty}^x f_x(x)$$

Given RVs  $x \in \mathcal{X}, y \in \mathcal{Y}$ , *marginalization* (sum rule) and *conditioning* (product rule) are given by:

$$f_x(x) = \sum_{y \in \mathcal{Y}} f_{xy}(x, y) \quad f_{x|y}(x|y) = \frac{f_{xy}(x, y)}{f_y(y)}$$

Which also apply to conditional PDFs; given  $z \in \mathcal{Z}$ :

$$f_{x|z}(x|z) = \sum_{y \in \mathcal{Y}} f_{xy|z}(x, y|z) \quad f_{x|y, z}(x|y, z) = \frac{f_{xy|z}(x, y|z)}{f_{y|z}(y|z)}$$

The above leads to *total probability theorem* and *Bayes' theorem*:

$$f_x(x) = \sum_{y \in \mathcal{Y}} f_{x|y}(x|y) f_y(y) \quad f_{x|y}(x|y) = \frac{f_{y|x}(y|x) f_x(x)}{f_y(y)}$$

$x$  and  $y$  are *independent* if  $f(x|y) = f(x) \iff f(x, y) = f(x)f(y)$ .

$x$  and  $y$  are *conditionally independent* given  $z$  if

$$f(x|y, z) = f(x|z) \iff f(x, y|z) = f(x|z)f(y|z)$$

The *expected value* or *mean* of  $\mathbf{x}$  is  $\boldsymbol{\mu}_x = \mathbb{E}[\mathbf{x}] = \sum_{\mathbf{x} \in \mathcal{X}} \mathbf{x} f_x(\mathbf{x})$

The (co)variance of  $\mathbf{x}$  is

$$\boldsymbol{\Sigma}_x = \text{Var}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^T$$

$\text{Var}[x] = \sigma_x^2 = \mathbb{E}[x^2] - (\mathbb{E}[x])^2$  in 1D.  $\boldsymbol{\Sigma}$  is symmetric positive (semi-)definite.

*Law of the Unconscious Statistician:*

$$\mathcal{Y} = \{ y \mid y = g(x), x \in \mathcal{X} \} \implies \mathbb{E}[y] = \mathbb{E}[g(x)]$$

**Change of Variables:** Given  $f_y(y)$ , let  $x = g(y)$ ; then

$$\text{Discrete: } f_x(x_j) = \sum_{y_j, i \in \mathcal{Y}_j} f_y(y_{j,i}) \quad \text{Continuous: } f_x(x) = \frac{1}{\frac{dg(y)}{dy}} f_y(y)$$

where  $\mathcal{Y}_j$  is the set of all  $y_{j,i}$  such that  $g(y_{j,i}) = x_j$  (discrete); and  $g(y)$  is required to be continuously differentiable and strictly increasing/decreasing (continuous only).

## 10 Bayesian Tracking

Formulation: computing  $f(\mathbf{x}_k | \mathbf{y}_{1:k})$  from  $f(\mathbf{x}_{k-1} | \mathbf{y}_{1:k-1})$ , given some motion model  $\mathbf{x}_k = \mathbf{f}_{k-1}(\mathbf{x}_{k-1}, \mathbf{v}_{k-1})$  and observation model  $\mathbf{y}_k = \mathbf{h}_k(\mathbf{x}_k, \mathbf{w}_k)$ , where  $\mathbf{x}_k, \mathbf{w}_k$  are independent with known PDFs, and assuming the *Markov property*:  $f(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{x}_{k-2}, \dots) = f(\mathbf{x}_k | \mathbf{x}_{k-1})$ .

**Prior Update:**  $f(\mathbf{x}_k | \mathbf{y}_{1:k-1}) = \sum_{\mathbf{x}_{k-1} \in \mathcal{X}} f(\mathbf{x}_k | \mathbf{x}_{k-1}) f(\mathbf{x}_{k-1} | \mathbf{y}_{1:k-1})$

$f(\mathbf{x}_k | \mathbf{x}_{k-1})$  can be computed from motion model and noise distribution by change of variables.

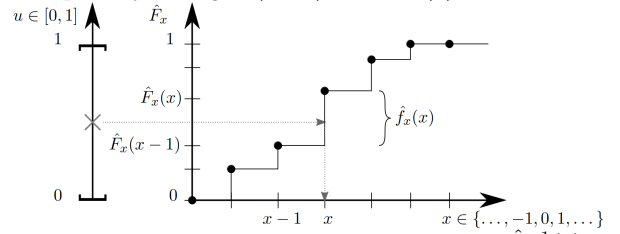
**Measurement Update:**  $f(\mathbf{x}_k | \mathbf{y}_{1:k}) = \frac{f(\mathbf{y}_k | \mathbf{x}_k) f(\mathbf{x}_k | \mathbf{y}_{1:k-1})}{\sum_{\mathcal{X}_k \in \mathcal{X}} f(\mathbf{y}_k | \mathcal{X}_k) f(\mathcal{X}_k | \mathbf{y}_{1:k-1})}$

$f(\mathbf{y}_k | \mathbf{x}_k)$  can be computed from measurement model and noise distribution by change of variables. The term in the denominator is for normalization.

## 11 Sampling From Distributions

Given  $u$  sampled from  $f_u(u)$ , a continuous uniform distribution over  $(0, 1)$ :

**One Variable, Discrete:** Given a desired PDF  $\hat{f}_x(x)$ , remap so  $\mathcal{X} = \mathbb{Z}$ , obtain a sample  $x$  by solving  $\hat{F}_x(x-1) < u \leq \hat{F}_x(x)$ .



**One Variable, Continuous:** Obtain sample by  $x = \hat{F}_x^{-1}(u)$  if CDF is invertible; otherwise pick any  $x$  value satisfying  $\hat{F}_x(x) = u$ .

**Multiple Variables (Finite Sample Space):** Given a desired discrete joint PDF  $\hat{f}_{xy}(x, y)$ , let  $\mathcal{Z} = 1, 2, \dots, N$  where  $N$  is the number of unique combinations of  $(x, y)$ ; define  $\hat{f}_z(z) = \hat{f}_{xy}(x, y)$  for some one-to-one mapping  $z \leftrightarrow (x, y)$ ; sample  $z$  using the standard method, then use the mapping to convert to  $(x, y)$ .

**Multiple Variables (General):** Find  $\hat{f}_y(y)$  and  $\hat{f}_{xy}(x, y) = \hat{f}_{x|y}(x|y)\hat{f}_y(y)$ ; sample  $\hat{f}_y(y)$  for  $y$ , substitute into  $\hat{f}_{x|y}(x|y)$  to sample for  $x$ .

## 12 Extracting Estimates from PDFs

Goal: Given observation  $\mathbf{y}$  and observation model  $f(\mathbf{y}|\mathbf{x})$ , estimate  $\mathbf{x}$ .

**Maximum Likelihood (ML):**  $\hat{\mathbf{x}}^{\text{ML}} = \arg \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{y}|\mathbf{x})$

Finds peak of distribution; best used when  $f(\mathbf{x})$  is not known. More sensitive to outliers and modelling error, since the distribution peak might not be the most robust option.

**Maximum a Posteriori (MAP):**  $\hat{\mathbf{x}}^{\text{MAP}} = \arg \max_{\mathbf{x}} f(\mathbf{y}|\mathbf{x})f(\mathbf{x})$

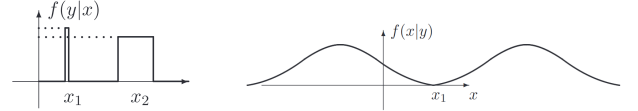
Finds peak of distribution  $f(\mathbf{x}|\mathbf{y})$  given known distribution of  $f(\mathbf{x})$ . Suffers from the same sensitivity issues of ML.

**Minimum Mean Squared Error (MMSE):**

$$\hat{\mathbf{x}}^{\text{MMSE}} = \arg \min_{\hat{\mathbf{x}}} \mathbb{E}_{\mathbf{x}|\mathbf{y}} \left[ (\mathbf{x} - \hat{\mathbf{x}})^T (\mathbf{x} - \hat{\mathbf{x}}) \right] = \mathbb{E}[\mathbf{x}|\mathbf{y}]$$

Finds expected value/mean of  $\mathbf{x}$  conditioned on  $\mathbf{y}$ . Typically more robust to outliers and modelling errors as the entire distribution is considered, but bad for bimodal distributions.

Undesirable cases shown below: Left: ML picks  $x_1$ , which is not robust. Right: MMSE picks  $x_1$ , which has zero probability.



## 13 Properties of Gaussian Distributions

1-dimensional:  $f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right)$

N-dimensional:  $f(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N \det \boldsymbol{\Sigma}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$

**Central Limit Theorem:** Let  $y_1, \dots, y_n$  be a sequence of  $n$  independent random variables drawn from (possibly different) distributions with finite population mean  $\mu$  and variance  $\sigma^2$  and let  $\bar{y} = y_1 + \dots + y_n$ , then for  $n \rightarrow \infty$ ,  $\bar{y} \sim \mathcal{N}(n\mu, n\sigma^2)$ .

A multivariate Gaussian  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  can be partitioned:

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \implies \mathbf{x} \sim f(\mathbf{x}_1, \mathbf{x}_2)$$

Then  $f(\mathbf{x}_1) = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$ ,  $f(\mathbf{x}_2) = \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$ , and

$$f(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N}(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21})$$

Gaussians are preserved in any linear transformation:

$$\mathbf{y}_1 \sim \mathcal{N}(\mathbf{a}_1, \mathbf{B}_1), \mathbf{y}_2 \sim \mathcal{N}(\mathbf{a}_2, \mathbf{B}_2)$$

$$\implies \mathbf{C}_1 \mathbf{y}_1 + \mathbf{C}_2 \mathbf{y}_2 \sim \mathcal{N}(\mathbf{C}_1 \mathbf{a}_1 + \mathbf{C}_2 \mathbf{a}_2, \mathbf{C}_1^T \mathbf{B}_1 \mathbf{C}_1 + \mathbf{C}_2^T \mathbf{B}_2 \mathbf{C}_2)$$

For nonlinear  $\mathbf{y} = \mathbf{g}(\mathbf{x})$ , if  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$ , then  $\mathbf{y} \sim \mathcal{N}(\mathbf{g}(\boldsymbol{\mu}_x), \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{A}^T)$  (approximately) where  $\mathbf{A}$  is the Jacobian of  $\mathbf{g}(\mathbf{x})$  at  $\boldsymbol{\mu}_x$ , if  $\mathbf{g}$  is roughly linear within  $\pm 3\sigma$  of  $\boldsymbol{\mu}_x$ .  $\mathbf{y} \sim \mathcal{N}(\mathbf{g}(\boldsymbol{\mu}_x), a^2 \sigma_x^2)$  in 1D where  $a = g'(\boldsymbol{\mu}_x)$ .

The product of multiple Gaussians is a Gaussian if normalized, with

$$\mathcal{N}(\boldsymbol{\mu}, \sigma^2) = \beta \prod_i \mathcal{N}(\mu_i, \sigma_i^2) \quad \frac{1}{\sigma^2} = \sum_i \frac{1}{\sigma_i^2} \quad \frac{\boldsymbol{\mu}}{\sigma^2} = \sum_i \frac{\mu_i}{\sigma_i^2}$$

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \beta \prod_i \mathcal{N}(\mu_i, \Sigma_i) \quad \boldsymbol{\Sigma}^{-1} = \sum_i \boldsymbol{\Sigma}_i^{-1} \quad \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} = \sum_i \boldsymbol{\Sigma}_i^{-1} \mu_i$$