

Data-Driven Bayesian and Non-Bayesian Parameter Estimation

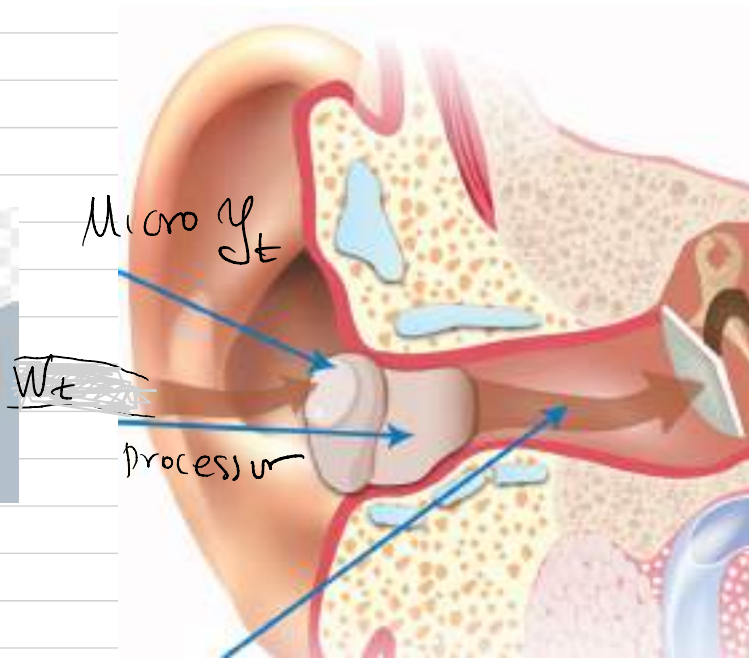
1



After detecting an airplane we would like to estimate its **position and speed** using the radar data $\{x_1, \dots, x_N\}$



Speaker

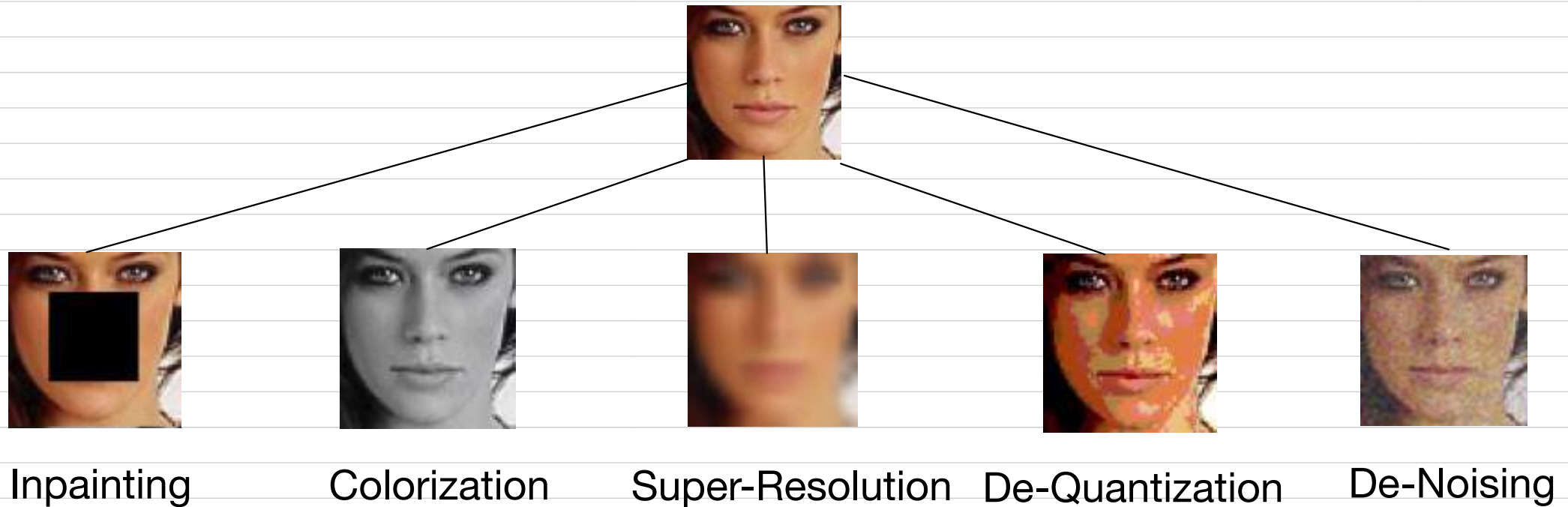


$$y_t = w_t + \underbrace{h_1 x_{t-\tau_1} + \dots + h_k x_{t-\tau_k}}_{\text{echo cancellation. (removal)}}$$

In many problems we would like to measure X but instead we measure Y . Can we recover (estimate) X by processing Y ?

X plays the role of “parameters”

Image Restoration (inverse problems)



Classical Parameter Estimation

Bayesian Approach

Data measurements $X = \{x_1, \dots, x_N\}$

X is realization of **random vector** \mathcal{X} with probability density $f(X|\theta)$ which known up to parameter vector θ

What of θ ?

θ is a realization of a **random vector** ϑ with density $p(\theta)$.
It expresses the **prior** knowledge about the parameters from observations made in past

Random vectors \mathcal{X}, ϑ are **statistically related** and relationship captured by the joint density

$$f(X, \theta) = f(X|\theta)p(\theta)$$

Problem of Interest

“Nature” selects a realization θ of ϑ that follows density $p(\theta)$

For the given θ “Nature” generates a realization X of \mathcal{X} that follows $f(X|\theta)$

From the measurements X and the prior information $f(X|\theta)$, $p(\theta)$

Estimate the θ that gave rise to the measurements

What is an “Estimator”??

Assume $\{x_1, \dots, x_N\}$ realizations of a random variable χ with mean μ .

We want to estimate μ

$$\hat{\mu}_1 = \frac{x_1 + \dots + x_N}{N} \qquad \hat{\mu}_2 = \frac{e^{x_1} + \dots + e^{x_N}}{N^2}$$

ANY function of the measurements can play the role of an estimator!!

Performance Computation

Define a cost function $C(\hat{\theta}, \theta)$

Compute **Average Cost**

$$C(\hat{\theta}) = \mathbb{E}_{x, \vartheta} [C(\hat{\theta}(X), \vartheta)]$$

Minimize Average Cost with respect to function $\hat{\theta}(X)$

$$G(U, X) = \int C(U, \theta) f(\theta|X) d\theta = \frac{\int C(U, \theta) f(X|\theta) p(\theta) d\theta}{\int f(X|\theta) p(\theta) d\theta}$$

$$\hat{\theta}_o(X) = \arg \min_U G(U, X)$$

Examples

$$C(U, \theta) = \|U - \theta\|^2 \quad \text{MMSE}$$

$$\hat{\theta}_{\text{MMSE}}(X) = \mathbb{E}[\vartheta|X] = \int \theta f(\theta|X) d\theta = \frac{\int \theta f(X|\theta) p(\theta) d\theta}{\int f(X|\theta) p(\theta) d\theta}$$

Minimum
Mean Square
Error

$$C(U, \theta) = \|U - \theta\|_{L_1} = |U_1 - \theta_1| + \dots + |U_k - \theta_k| \quad \text{MAE}$$

for scalar U, θ

$$\hat{\theta}_{\text{MAE}}(X) = \arg \left\{ U : \int_{-\infty}^U f(\theta|X) d\theta = \frac{1}{2} \right\}$$

Minimum
Mean Absolute
Error

$$C(U, \theta) = \begin{cases} 1 & \text{when } \|U - \theta\| \geq \delta, \\ 0 & \text{when } \|U - \theta\| < \delta, \end{cases} \quad \delta \rightarrow 0 \quad \text{MAP}$$

$$\hat{\theta}_{\text{MAP}}(X) = \arg \max_{\theta} f(\theta|X) = \arg \max_{\theta} f(X|\theta) p(\theta)$$

Maximum
Aposteriori
Probability

Basic Tools



Neural Networks

Special class of **parametric** functions $u(X, \alpha)$ where α the network parameters

Any function $v(X)$ can be approximated **arbitrarily close** by a neural network of sufficiently high order

Searching over α in a neural network $u(X, \alpha)$ corresponds to searching over any function $v(X)$ when the size of the network becomes arbitrarily large

Law of Large Numbers (LLN)

\mathcal{X} random and $\{X_1, X_2, \dots, X_N\}$ realizations

Let $G(X)$ be a deterministic function, then

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N G(X_i) = \mathbb{E}_{\mathcal{X}} [G(\mathcal{X})] = \int G(X) f(X) dX$$

Gradient Descent

Deterministic function $J(\theta)$ and interested in $\min_{\theta} J(\theta)$

We can apply: $\theta_t = \theta_{t-1} - \mu \nabla_{\theta} J(\theta_{t-1})$

Stochastic Gradient Descent

$$J(\theta) = \mathbb{E}_{\mathcal{X}} [G(\mathcal{X}, \theta)]$$

Instead of $f(X)$ we have $\{X_1, \dots, X_N\}$ then

$$\theta_t = \theta_{t-1} - \mu \nabla_{\theta} G(X_t, \theta_{t-1}), \quad \mu > 0$$

Data - Driven Version

9

Classical estimation assumes availability of $f(X|\theta), p(\theta)$

Equivalent to joint density $f(X, \theta)$

$f(X, \theta)$ expresses the random relationship between \mathcal{X} and ϑ

If $f(X, \theta)$ to be replaced by data then we need **pairs**

$$\{(X_1, \theta_1), (X_2, \theta_2), \dots, (X_N, \theta_N)\}$$

The general estimator function $\hat{\theta}(X)$ is replaced by $u(X, \alpha)$
a neural network with parameters α

The optimization becomes

$$\min_{\hat{\theta}(X)} \mathbb{E}_{\mathcal{X}, \vartheta} [C(\hat{\theta}(\mathcal{X}), \vartheta)] \Rightarrow \min_{\alpha} \mathbb{E}_{\mathcal{X}, \vartheta} [C(u(\mathcal{X}, \alpha), \vartheta)]$$



Apply Stochastic Gradient Descent

$$\alpha_t = \alpha_{t-1} - \mu [\mathbb{J}_\alpha u(X_t, \alpha_{t-1})]^\top \nabla_U \mathbf{C}(u(X_t, \alpha_{t-1}), \theta_t)$$

Alternatively replace expectation using LLN

$$\min_{\alpha} \sum_{i=1}^N \mathbf{C}(u(X_i, \alpha), \theta_i)$$

Solve using Gradient Descent with respect to α

If limit is α_o then we expect

$$u(X, \alpha_o) \approx \hat{\theta}_o(X)$$

Non-Bayesian Estimation

“Nature” selects θ and generates measurements X that follow $f(X|\theta)$

We know $f(X|\theta)$ up to a set of parameters θ

Parameters θ are deterministic and unknown

Problem: Given measurement vector X estimate θ

$$\hat{\theta}_{\text{MLE}}(X) = \arg \max_{\theta} f(X|\theta)$$

Optimality? (Asymptotic)

$$\mathbb{E}[\|\hat{\theta}(X) - \theta\|^2] \geq \text{CRLB} \quad \frac{\mathbb{E}[\|\hat{\theta}_{\text{MLE}}(X) - \theta\|^2]}{\text{CRLB}} \rightarrow 1$$

(CRAMER-RAO Lower Bound)
 $|X| \rightarrow \infty$

Data - Driven Version

In the classical setup, parameter estimation makes sense only if there is parametric density function $f(X|\theta)$

We will sacrifice generality in order to define a meaningful parameter estimation problem that can be formulated under a data-driven setup

We will define $f(X|\theta)$ **indirectly!!**

We start with random vector $Z \sim g(Z)$

We consider a deterministic transformation $T(Z, \theta)$ which is known up to some parameter vector θ

We define the random vector $X = T(Z, \theta)$ which has density $X \sim f(X|\theta)$

Assume instead of $g(Z)$ we have dataset $\{Z_1, \dots, Z_m\}$ with independent realizations of Z

Assume that we are given a dataset $\{X_1, \dots, X_n\}$ with independent realizations of X all following $f(X|\theta)$ with the same θ .

Problem:

Using $\{Z_1, \dots, Z_m\}$ as a representative of the density $g(Z)$

Assuming knowledge of the transformation $T(Z, \theta)$ up to the unknown parameters θ

For every collection of data $\{X_1, \dots, X_n\}$ **estimate the θ** that has generated them.

IF for the data we had $X_i = T(Z_i, \theta)$ the **problem would have been simple.**

We could form some distance between the X s and the $T(Z, \theta)$ s and minimize over θ

BUT the two datasets $\{Z_1, \dots, Z_m\}$ and $\{X_1, \dots, X_n\}$ are **independent and unrelated.**

Moment Matching

If $X = T(Z, \theta)$ then

$$\mathbb{E}[X^p] = \mathbb{E}[(T(Z, \theta))^p]$$

$$\frac{1}{n} \sum_{i=1}^n (X_i)^p = \frac{1}{m} \sum_{j=1}^m (T(Z_j, \theta))^p$$

$$P = p_1, p_2, \dots$$

Sufficient # of equations to solve for θ .

Many different choices for moments

Moment estimates are **Notoriously NON-ROBUST**

Density Matching

We would like to find θ so that \mathcal{X} and $\mathbb{T}(Z, \theta)$ exhibit **the same statistical behavior**

We would like to find θ so that $\{X_1, \dots, X_n\}$ and $\{\mathbb{T}(Z_1, \theta), \dots, \mathbb{T}(Z_m, \theta)\}$ exhibit the same statistical behavior

There exists an interesting methodology developed for Generative Modeling and our problem constitutes a special case.

Generative Adversarial Networks (GANs)

The random vector \mathcal{X} follows $f(X)$ and the random vector \mathcal{Z} follows $g(Z)$

We would like to find a transformation (**generator**) $G(Z)$ such that $\mathcal{Y} = G(\mathcal{Z})$ follows $f(X)$

To solve the problem we are going to design a second function $D(X)$ (**discriminator**) by considering the **adversarial problem**

$$\min_{G(Z)} \max_{D(X)} J(D, G)$$

2014 Goodfellow et al.

where

$$J(D, G) = \mathbb{E}_x [\log D(x)] + \mathbb{E}_z [\log (1 - D(G(z)))]$$

→ $G(z)$ such that $\mathcal{Y} = G(\mathcal{Z}) \sim f(x)$

Extensions

$$J(D, G) = \mathbb{E}_x [\phi(D(X))] + \mathbb{E}_z [\psi(D(G(Z)))]$$

$$\min_{G(Z)} \max_{D(X)} J(D, G)$$

$\psi'(z) < 0$
 $\omega'(r) > 0$
 $\phi'(z) = -\bar{\omega}'(z) \psi'(z)$

designs the correct “generator”

Data-Driven Setup

Under a data-driven setup $G(Z) \rightarrow G(Z, \theta)$ and $D(X) \rightarrow D(X, \vartheta)$ with the generator and discriminator becoming parametric transformations

This is exactly the same with our problem where the “generator” is $T(Z, \theta)$.

$$J(\theta, \vartheta) = \frac{1}{n} \sum_{i=1}^n \phi(D(X_i, \vartheta)) + \frac{1}{m} \sum_{j=1}^m \psi(D(G(Z_j, \theta), \vartheta))$$

$$\min_{\theta} \max_{\vartheta} J(\theta, \vartheta)$$

Optimum $\theta_o \Rightarrow G(Z, \theta_o)$

If Z realization of \mathcal{Z} following $g(Z)$ then $G(Z, \theta_o)$ realization of \mathcal{X} following $f(X)$.

EXAMPLE

If $\{X_1, X_2, \dots, X_m\}$ (many) human faces and
 $\{z_1, z_2, \dots, z_m\}$ (many) Gaussian random vectors

Can design function $G(Z, \theta_o)$ such that when Z
 Gaussian vector $G(Z, \theta_o)$ is **A HUMAN FACE**

DATASET CELIBA (x_1, x_2, \dots)



True faces

Generate z_1, z_2, \dots Gaussian vectors
↓
 $G(z_1, \theta_0)$ $G(z_2, \theta_0)$ \dots



Synthetic faces

Assume a random vector \mathcal{X} is described by a **Generative Model**: $g(Z), G(Z)$ instead of a density $f(X)$

Advantage: If we can easily generate realizations of $g(Z)$ then we transform them and generate realizations of $f(X)$

Inverse Problems

For \mathcal{X} we have generative model $g(Z), G(Z)$

If X , a realization of \mathcal{X} , undergoes a transformation $Y = F(X)$ then X can be recovered from Y by first recovering the Z that generates X as $X = G(Z)$.



X



Y

$$X = G(Z) \quad Z \sim g(Z)$$

Gaussian.

from Y estimate \hat{Z}

$$\text{then } \underline{\hat{X} = G(\hat{Z})}$$



De-Quantization

Colorization

X^1 X^2 Y $\hat{X}^1 = G_1(\hat{z}^1)$ $\hat{X}^2 = G_2(\hat{z}^2)$

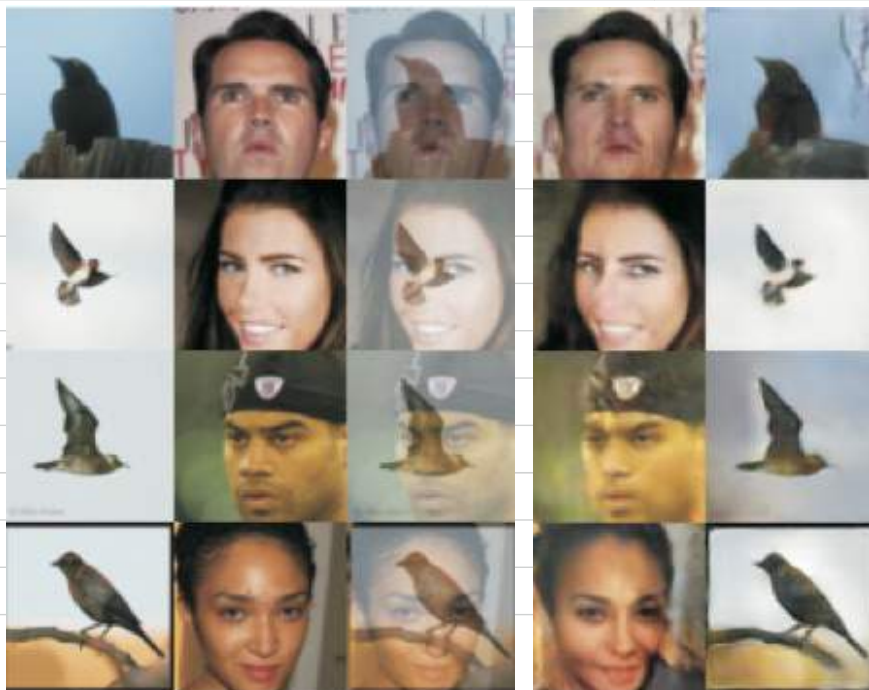


Image Separation