

(A Quick) Probability Review

Reading:

- Go over handouts 2–5 in EE 420x notes.

Basic probability rules:

(1) $P[\Omega] = 1, P[\emptyset] = 0, 0 \leq P[A] \leq 1;$
 $P[\cup_{i=1}^{\infty} A_i] = \sum_{i=1}^{\infty} P[A_i]$ if $\underbrace{A_i \cap A_j}_{A_i \text{ and } A_j \text{ disjoint}} = \emptyset$ for all $i \neq j;$

(2) $P[A \cup B] = P[A] + P[B] - P[A \cap B], P[A^c] = 1 - P[A];$

(3) If $A \perp B$, then $P[A \cap B] = P[A] \cdot P[B];$

(4)

$$P[A | B] = \frac{P[A \cap B]}{P[B]} \quad (\text{conditional probability})$$

or

$$P[A \cap B] = P[A | B] \cdot P[B] \quad (\text{chain rule});$$

(5)

$P[A] = P[A | B_1] P[B_1] + \dots + P[A | B_n] P[B_n]$
if B_1, B_2, \dots, B_n form a *partition* of Ω ;

(6) Bayes rule:

$$P[A | B] = \frac{P[B | A] P[A]}{P[B]};$$

(7)

$$\begin{aligned} E[aX + bY + c] &= a \cdot E[X] + b \cdot E[Y] + c \\ \text{var}(aX + bY + c) &= a^2 \text{var}(X) + b^2 \text{var}(Y) \\ &\quad + 2ab \cdot \text{cov}(X, Y) \end{aligned}$$

where a, b , and c are constants and X and Y are random variables.

(7') A vector/matrix version of (7):

$$\begin{aligned} E[A\mathbf{X} + B\mathbf{Y} + \mathbf{c}] &= A E[\mathbf{X}] + B E[\mathbf{Y}] + \mathbf{c} \\ \text{cov}(A\mathbf{X} + B\mathbf{Y} + \mathbf{c}) &= A \text{cov}(\mathbf{X}) A^T + B \text{cov}(\mathbf{Y}) B^T \\ &\quad + A \text{cov}(\mathbf{X}, \mathbf{Y}) B^T + B \text{cov}(\mathbf{Y}, \mathbf{X}) A^T \end{aligned}$$

where “ T ” denotes a transpose and

$$\text{cov}(\mathbf{X}, \mathbf{Y}) = E\{(\mathbf{X} - E[\mathbf{X}])(\mathbf{Y} - E[\mathbf{Y}])^T\}.$$

(To refresh memory about covariance and its properties, see p. 12 of handout 5 in EE 420x notes. For random vectors, see handout 7 in EE 420x notes, particularly pp. 1–15.)

Some useful theorems:

(1) (handout 5 in EE 420x notes)

$$\begin{aligned}E[X] &= E_Y[E_{X|Y}[X|Y]] \\E[g(X) \cdot h(Y) | Y = y] &= h(y) \cdot E[g(X) | Y = y] \\E[g(X) \cdot h(Y)] &= E[h(Y) \cdot E[g(X) | Y]];\end{aligned}$$

The vector version of (1) is the same — just put bold letters.

(2)

$$\text{var}(X) = E[\text{var}(X | Y)] + \text{var}(E[X|Y]);$$

The vector/matrix version of (2) is:

(2')

$$\underbrace{\text{cov}(\mathbf{X})}_{\text{variance/covariance matrix of } \mathbf{X}} = E[\text{cov}(\mathbf{X} | \mathbf{Y})] + \text{cov}(E[\mathbf{X} | \mathbf{Y}]);$$

(3)

$$\text{cov}(X, Y) = \text{E}[\text{cov}(X, Y | Z)] + \text{cov}(\text{E}[X|Z], \text{E}[Y|Z]).$$

(4) Transformation:

$$\mathbf{Y} = \mathbf{g}(\mathbf{X}) \iff \begin{array}{l} Y_1 = g_1(X_1, \dots, X_n) \\ \vdots \\ Y_n = g_n(X_1, \dots, X_n) \end{array}$$

then

$$p_Y(\mathbf{y}) = p_X(h_1(y_1), \dots, h_n(y_n)) \cdot |J|$$

where $h(\cdot)$ is the unique inverse of $g(\cdot)$ and

$$J = \left| \frac{\partial \mathbf{x}}{\partial \mathbf{y}^T} \right| = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \cdots & \frac{\partial x_1}{\partial y_n} \\ \vdots & \vdots & \vdots \\ \frac{\partial x_n}{\partial y_1} & \cdots & \frac{\partial x_n}{\partial y_n} \end{vmatrix}$$

Print and read the handout “Probability distributions” from the Supplementary material section on WebCT. Bring it with you to the midterm exam.

Estimator Performance

We now continue with the DC-level estimation example from handout # 0.

This handout describes the “classical world” and, therefore, arguments are made in this context. Bayesian arguments will be somewhat different.

Consider the following two estimators:

$$\hat{A}_1 = \frac{1}{N} \sum_{n=0}^{N-1} x[n]$$
$$\hat{A}_2 = x[0].$$

Note that \hat{A}_1 is the ML estimate of A — it maximizes the likelihood function, see handout # 0. (Interestingly, the ML estimate of A is the same regardless of whether σ^2 is known or not.) It is also intuitively appealing — A is the *average* level of $x[n]$ (since $w[n]$ is zero mean).

Which estimator is better?

For a given realization, it is possible that either \hat{A}_1 or \hat{A}_2 is closer to A . Hence, we need statistical analysis to answer this question.

Estimator Performance (cont.)

Substitute the measurement model to perform statistical analysis. We have

$$\hat{A}_1 = \frac{1}{N} \sum_{n=0}^{N-1} \{A + w[n]\}$$

$$\hat{A}_2 = A + w[0].$$

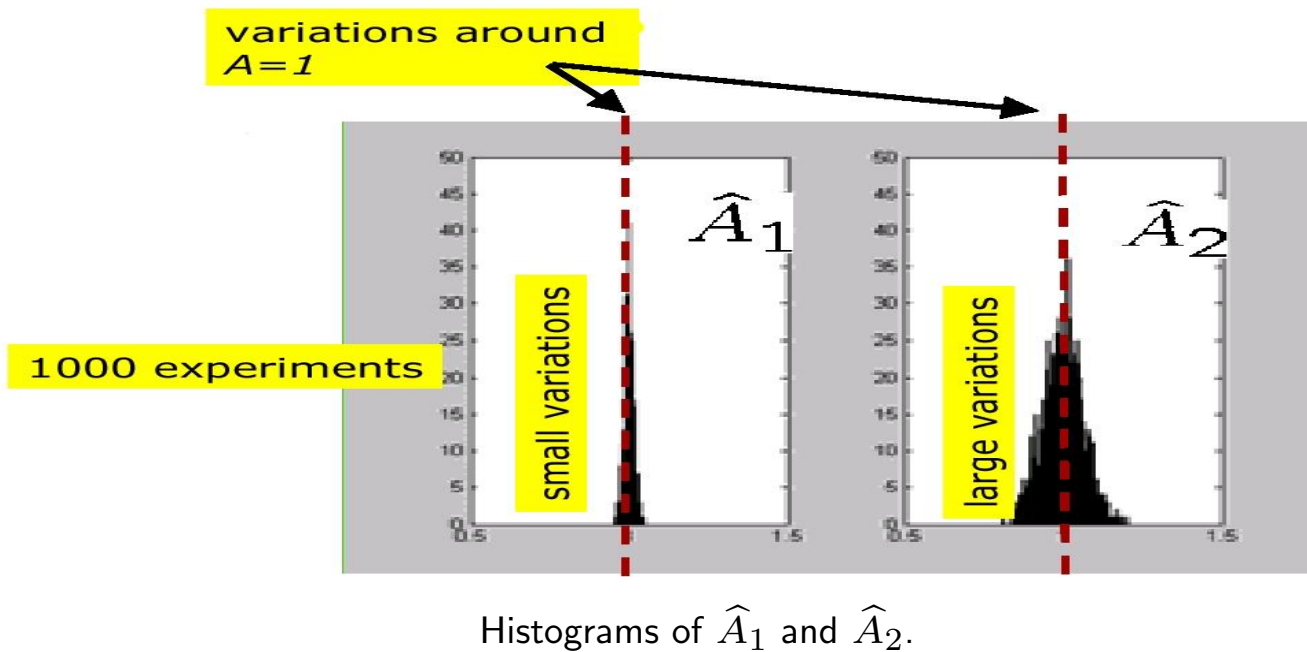
Take expectation:

$$\mathbb{E}_x[\hat{A}_1] = \frac{1}{N} \sum_{n=0}^{N-1} (A + \overbrace{\mathbb{E}\{w[n]\}}^0) = A$$

$$\mathbb{E}_x[\hat{A}_2] = A + \overbrace{\mathbb{E}\{w[n]\}}^0 = A.$$

On average, both estimators are around the correct value (i.e. they are unbiased).

Estimator Performance (cont.)



But \hat{A}_1 is better than \hat{A}_2 because its pdf is more concentrated around the true value.

“On average, \hat{A}_1 is closer to $A = 1$.”

Proof.

$$E_x[\hat{A}_1] = E[\hat{A}_2] = A$$

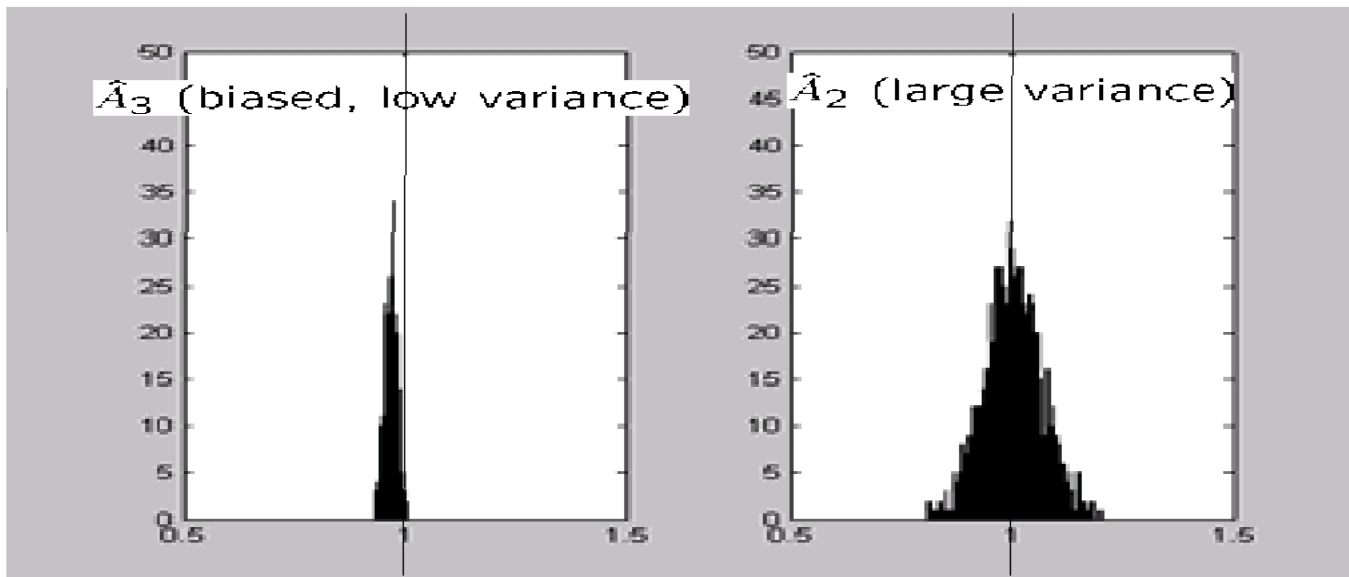
$$\text{var}_x(\hat{A}_1) = E[(\hat{A}_1 - \overbrace{E[\hat{A}_1]}^A)^2] = \frac{1}{N^2} \sum_{n=0}^{N-1} \text{var}_x(w[n]) = \frac{\sigma^2}{N}$$

$$\text{var}_x(\hat{A}_2) = \sigma^2.$$

□

But, what is the justification for taking these expectations? We implicitly assume that we could repeat this experiment many times and plot the histogram (say) of the resulting estimates. This is what Bayesians criticize saying that, in the classical approach, “data that have never been observed are used for inference.” (As you can guess, Bayesians do not need this virtual-data argument.) Suppose we are fine with the classical argument and let us continue.

Which Estimator is the Best?



Notation. Bias and mean-square error (MSE) of estimator $\hat{\theta}$:

$$b(\theta) = E_x[\hat{\theta}] - \theta$$

$$\text{MSE}(\hat{\theta}) = E_x[(\hat{\theta} - \theta)^2] = \text{var}(\hat{\theta}) + b(\theta)^2.$$

We wish to minimize the MSE, which leads to an MMSE estimator. MSE expression:

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= E_x[(\hat{\theta} - E_x[\hat{\theta}] + E_x[\hat{\theta}] - \theta)^2] \\ &= \underbrace{E_x[(\hat{\theta} - E_x[\hat{\theta}])^2]}_{\text{var}(\hat{\theta})} + \underbrace{(E_x[\hat{\theta}] - \theta)^2}_{b(\theta)^2} \\ &\quad + \underbrace{2 E_x[(\hat{\theta} - E_x[\hat{\theta}]) \cdot (E_x[\hat{\theta}] - \theta)]}_{0} \end{aligned}$$

Note that the above $\text{MSE}(\hat{\theta})$ is different from Bayesian MSE. Since in Bayesian inference we assign a prior distribution to θ , Bayesian MSE would be obtained by taking the expectation of $\text{MSE}(\hat{\theta})$ with respect to θ .

In the classical setup, minimizing the MSE is our key objective, but the use of this criterion leads to unrealizable estimators $\implies \text{MSE}(\hat{\theta}) = \text{var}(\hat{\theta}) + b(\theta)^2$ is a strong function of θ and minimizing it over some family of estimators will usually produce an “optimal” estimator $\hat{\theta}$ that depends on θ .

Example 1. DC level in **Additive White Gaussian Noise (AWGN)**, see also Section 2.4 in Kay-I:

$$x[n] = A + \underbrace{w[n]}_{\text{AWGN}}$$

$$w[n] \sim \mathcal{N}(0, \sigma^2), \quad n = 0, 1, \dots, N - 1.$$

Consider the following family of estimators of A :

$$\check{A} = a \bar{x}$$

where

$$\bar{x} = (1/N) \sum_{n=0}^{N-1} x[n] \quad (\text{sample mean}).$$

Here

$$\text{E}[\check{A}] = aA, \quad \text{var}(\check{A}) = a^2 \sigma^2 / N.$$

Find the best a (that minimizes the MSE for the given family).
 In other words, can we improve upon the sample mean?

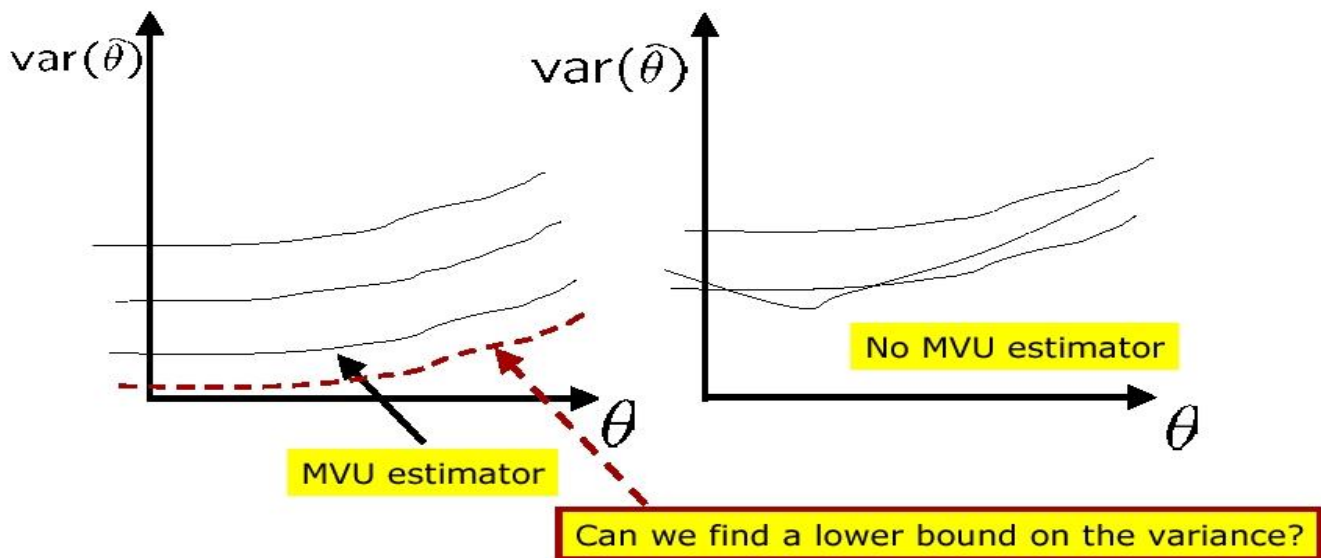
$$\begin{aligned} \text{MSE}(\check{A}) &= a^2\sigma^2/N + \overbrace{(aA - A)^2}^{\text{depends on } A} \\ \frac{d\text{MSE}}{da} &= 2a\sigma^2/N + 2(aA - A)A = 0 \\ a_{\text{opt}} &= \frac{A^2}{A^2 + \sigma^2/N} \end{aligned}$$

depends on the unknown parameter A . Hence not useful (at least not directly). Observe the “shrinkage” form of the above “estimator.”

Minimum Variance Unbiased (MVU) Estimation

How do we construct a “realizable” estimator?

An idea: Constrain the bias to be zero and then minimize the estimator variance (which is equal to MSE in this scenario since the bias is zero) for all values of $\theta \implies$ MVU estimator.



MVU estimator does not always exist, as $\hat{\theta}$ must have smallest variance for all values of θ .¹

¹To emphasize the fact that the MVU estimator must have the smallest variance for all values of θ , B & D refer to it as *uniformly minimum variance unbiased* (UMVU).

Comments:

- Even if it exists for a particular problem, MVU estimator is not optimal in terms of minimizing the MSE and we may be able to do better.
- Unbiasedness is nice, but not the most important \implies we can relax this condition and consider biased estimators as well, e.g. by making them *asymptotically unbiased*. By relaxing the unbiasedness condition, it is possible to outperform the MVU estimators in terms of MSE, as shown in the following example. What we really care about is minimizing the MSE!

Example 2². Consider now estimating the variance σ^2 of independent, identically distributed (i.i.d.) zero-mean Gaussian observations, using the following estimator:

$$\hat{\sigma}^2 = a \cdot \frac{1}{N} \sum_{n=0}^{N-1} x^2[n] \quad (1)$$

where $a > 0$ is variable. If we choose $a = 1$, $\hat{\sigma}^2|_{a=1}$ will be unbiased³ with

$$\hat{\sigma}^2|_{a=1} = \hat{\sigma}_{\text{MVU}}^2 = \frac{1}{N} \sum_{n=0}^{N-1} x^2[n]. \quad (2)$$

²See also P. Stoica and R. Moses, "On biased estimators and the unbiased Cramér-Rao lower bound," *Signal Processing*, vol. 21, pp. 349–350, 1991.

³We will show later that this choice yields an MVU estimate.

Now, in general,

$$\mathbb{E}[\hat{\sigma}^2] = a \sigma^2$$

and

$$\begin{aligned} \text{MSE}(\hat{\sigma}^2) &= \mathbb{E}[(\hat{\sigma}^2 - \sigma^2)^2] \\ &= \mathbb{E}[\hat{\sigma}^4] + \sigma^4 - 2\sigma^2 \mathbb{E}[\hat{\sigma}^2] \\ &= \mathbb{E}[\hat{\sigma}^4] + \sigma^4(1 - 2a) \\ &= \frac{a^2}{N^2} \sum_{n_1=0}^{N-1} \sum_{n_2=0}^{N-1} \mathbb{E}\{x^2[n_1]x^2[n_2]\} \\ &\quad + \sigma^4(1 - 2a) \\ &= \frac{a^2}{N^2} [(N^2 - N)\sigma^4 + N \cdot \underbrace{\mathbb{E}\{x^4[n]\}}_{3\sigma^4}] + \sigma^4(1 - 2a) \\ &= \sigma^4 \cdot \left[a^2 \left(1 + \frac{2}{N}\right) + (1 - 2a) \right]. \end{aligned} \quad (3)$$

To evaluate the above expression, we have used the following facts:

- For $n_1 \neq n_2$, $\mathbb{E}\{x^2[n_1]x^2[n_2]\} = \mathbb{E}\{x^2[n_1]\} \cdot \mathbb{E}\{x^2[n_2]\} = \sigma^2 \cdot \sigma^2 = \sigma^4$.
- For $n_1 = n_2$, $\mathbb{E}\{x^2[n_1]x^2[n_2]\} = \mathbb{E}\{x^4[n_1]\} = 3\sigma^4$ (which is the fourth-order moment of a Gaussian distribution).

It can be easily shown that (3) is minimized for

$$a_{\text{OPT}} = \frac{N}{N+2}$$

yielding the estimator

$$\hat{\sigma}_{\star}^2 = a_{\text{OPT}} \cdot \frac{1}{N} \sum_{n=0}^{N-1} x^2[n]$$

whose MSE

$$\text{MSE}_{\text{MIN}} = \frac{2\sigma^4}{N+2}.$$

is minimum for the family of estimators in (1).

Comments:

- $\hat{\sigma}_{\star}^2$ is *biased* and has *smaller MSE* than the MVU estimator in (2):

$$\text{MSE}_{\text{MIN}} < \text{MSE}(\hat{\sigma}^2) \Big|_{a=1} = \frac{2\sigma^4}{N}.$$

- Note that we are able to construct a realizable estimator in this case — compare with Example 1 in this handout.
- For large N , $\hat{\sigma}_{\star}^2$ and $\hat{\sigma}_{\text{MVU}}^2$ are approximately the same since $N/(N+2) \rightarrow 1$ as $N \rightarrow \infty$. This also implies that $\hat{\sigma}_{\star}^2$ is *asymptotically unbiased*.

Note: I do not wish to completely dismiss bias considerations. For example, we may have two estimators $\hat{\theta}_1$ and $\hat{\theta}_2$ with

$$[\text{bias}(\hat{\theta}_1)]^2 \ll \text{var}(\hat{\theta}_1) \quad \text{and} \quad [\text{bias}(\hat{\theta}_2)]^2 \ll \text{var}(\hat{\theta}_2)$$

and

$$\text{MSE}(\hat{\theta}_1) \approx \text{MSE}(\hat{\theta}_2).$$

So, these two estimators are “equally good” as far as MSE is concerned. But, we may have

$$|\text{bias}(\hat{\theta}_1)| \ll |\text{bias}(\hat{\theta}_2)|$$

making $\hat{\theta}_1$ “more desirable” than $\hat{\theta}_2$. *Bias correction* methods have been developed for constructing estimators that have small bias. Hence, having small bias is typically a second-tier concern (compared with minimizing the MSE), but a valid one, particularly in the scenario outlined in this comment.

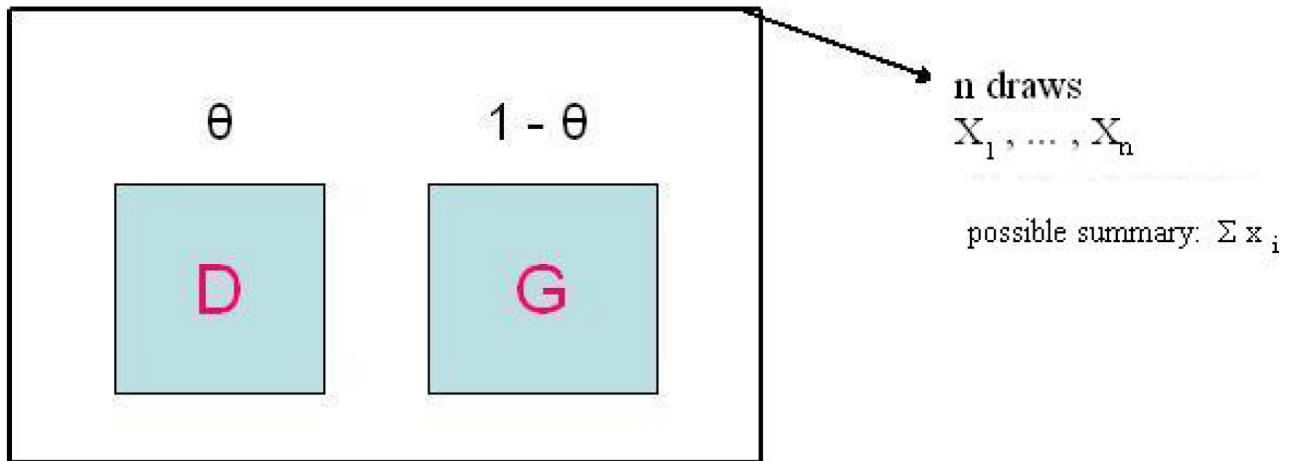
Sufficiency

Reading:

- Kay-I, chs. 5.3–5.4.

A function $T(\mathbf{x})$ of the observations \mathbf{x} *only* is called a *statistic*.

Example: A machine produces n items in succession with probability θ of producing a defective product. Suppose that there is no dependence in quality of the produced items.



Then, our statistical model is

$$p(\mathbf{x} | \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i}.$$

Is there any loss of information by keeping and recording only

$\sum_{i=1}^n x_i$? Answer:

$\left\{ \begin{array}{l} \text{Yes, we are dropping a lot of information. But} \\ \text{No, in terms of inference about } \theta \\ \text{(provided that the model is correct, of course).} \end{array} \right.$

We typically wish to separate out any aspects of the data that are irrelevant in the context of our model. In other words, we would like to reduce the data and deal only with the statistics “whose use involves no loss of information.” For example, we could save memory and store only the reduced data. What we mean by “no loss of information” is quantified in the following definition.

Definition. $T = T(\mathbf{x})$ is a sufficient statistic for θ if the conditional distribution of \mathbf{X} given $T(\mathbf{X}) = T$ *does not* involve θ :

$$p(\mathbf{x} | T(\mathbf{x}) = t; \theta) = p(\mathbf{x} | T(\mathbf{x}) = t).$$

Think of sufficient statistics as not throwing away any useful information about θ .

Note: if X is a random variable (RV), then $T = T(X)$ is a RV.

Trivial example: $T(\mathbf{x}) = \mathbf{x} \implies$ full data is always sufficient.

Example: Let us continue with the previous example. Here, $\mathbf{X} = [X_1, \dots, X_n]^T$ is the record of n Bernoulli trials with

probability θ , which can be written as

$$P\{X_i = x_i\} = \theta^{x_i}(1 - \theta)^{1-x_i}.$$

where x_i is 1 (defective) or 0 (not defective). Thus

$$P\{\mathbf{X} = \mathbf{x}\} = P\{X_1 = x_1, \dots, X_n = x_n\} = \theta^t(1 - \theta)^{n-t}$$

where $t = \sum_{i=1}^n x_i$. Now, $T = \sum_{i=1}^n X_i$ has a binomial distribution $\text{Bin}(n, \theta)$ and

$$\begin{aligned} p(\mathbf{x}|T(\mathbf{x}) = t; \theta) &= P\{\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = t\} \\ &= \frac{P\{\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = t\}}{\binom{n}{t} \theta^t (1 - \theta)^{n-t}} \\ &= \begin{cases} 0, & \text{if } \sum_{i=1}^n x_i \neq t \\ \frac{\theta^t (1 - \theta)^{n-t}}{\binom{n}{t} \theta^t (1 - \theta)^{n-t}} = \frac{1}{\binom{n}{t}}, & \text{otherwise} \end{cases} \end{aligned}$$

which is clearly not a function of θ . Here, we have used the (general) fact that

$$\{\mathbf{X} = \mathbf{x}\} \subset \{T(\mathbf{X}) = T(\mathbf{x})\}. \quad (4)$$

Thus, $T(\mathbf{x}) = \sum_{i=1}^n x_i$ is a sufficient statistic for θ .

In general, directly checking sufficiency is difficult because we need to compute the conditional distribution. Fortunately, we have the following theorem whose conditions are easy to verify.

Theorem. (Factorization Theorem) *A statistic $T(\mathbf{X})$ is sufficient for θ if and only if there exists a function $g(t, \theta)$ and a function $h(\mathbf{x})$ such that*

$$p(\mathbf{x}; \theta) = \underbrace{g(T(\mathbf{x}), \theta)}_{\substack{\text{parameters} \\ \text{coupled with} \\ \text{suff. stat.}}} \cdot h(\mathbf{x}).$$

Note: $T(\mathbf{x})$ must be a statistic, a function of data \mathbf{x} *only*.

Proof. To illustrate the idea of the proof and for simplicity, we concentrate on the discrete case. Suppose that $T(\mathbf{X})$ is sufficient. Then

$$\begin{aligned} p(\mathbf{x}; \theta) &= \overbrace{P\{\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = T(\mathbf{x})\}}^{P\{\mathbf{X} = \mathbf{x}\}, \text{ see (4)}} \\ &= P\{T(\mathbf{X}) = T(\mathbf{x})\} \underbrace{P\{\mathbf{X} = \mathbf{x} \mid T(\mathbf{X}) = T(\mathbf{x})\}}_{\substack{h(\mathbf{x}), \\ \text{by sufficiency}}} \\ &= g(T(\mathbf{x}), \theta) h(\mathbf{x}). \end{aligned}$$

Conversely,

$$\begin{aligned}
 P\{\mathbf{X} = \mathbf{x} \mid T(\mathbf{X}) = T(\mathbf{x})\} &= \frac{P\{\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = T(\mathbf{x})\}}{P\{T(\mathbf{X}) = T(\mathbf{x})\}} \\
 &= \frac{\overbrace{P\{\mathbf{X} = \mathbf{x}\}}^{p(\mathbf{x}; \theta)}}{\underbrace{P\{T(\mathbf{X}) = T(\mathbf{x})\}}_{\sum_{\mathbf{y}: T(\mathbf{y})=T(\mathbf{x})} p(\mathbf{y}; \theta)}} \\
 &= \frac{\overbrace{g(T(\mathbf{x}), \theta) h(\mathbf{x})}^{\text{by the assumption}}}{\sum_{\mathbf{y}: T(\mathbf{y})=T(\mathbf{x})} \underbrace{g(T(\mathbf{y}), \theta) h(\mathbf{y})}_{\text{by the assumption}}} \\
 &= \frac{g(T(\mathbf{x}), \theta) h(\mathbf{x})}{g(T(\mathbf{x}), \theta) \sum_{\mathbf{y}: T(\mathbf{y})=T(\mathbf{x})} h(\mathbf{y})} \\
 &= \frac{h(\mathbf{x})}{\sum_{\mathbf{y}: T(\mathbf{y})=T(\mathbf{x})} h(\mathbf{y})}
 \end{aligned}$$

which is *not* a function of θ . \square

Example: Suppose x_1, x_2, \dots, x_n , are i.i.d. $\mathcal{N}(\mu, \sigma^2)$. Let

$\boldsymbol{\theta} = [\mu, \sigma^2]^T$. Then

$$\begin{aligned} p(x_1, \dots, x_n; \boldsymbol{\theta}) &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{n\mu^2}{2\sigma^2} \right\} \exp \left\{ -\frac{1}{2\sigma^2} \left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i \right) \right\}. \end{aligned}$$

Clearly, $p(x_1, \dots, x_n; \boldsymbol{\theta})$ is itself a function of $\sum_{i=1}^n x_i$, $\sum_{i=1}^n x_i^2$, and $\boldsymbol{\theta}$ only and, upon applying the factorization theorem, we conclude that

$$\mathbf{T}(\mathbf{x}) = \mathbf{T}(x_1, x_2, \dots, x_n) = \begin{bmatrix} \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i^2 \end{bmatrix}$$

is sufficient for $\boldsymbol{\theta}$. Here, $h(\mathbf{x})$ is trivial: $h(\mathbf{x}) = 1$.

An equivalent (frequently used) sufficient statistic is

$$\mathbf{T}(\mathbf{x}) = \mathbf{T}(x_1, x_2, \dots, x_n) = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n x_i \\ \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \end{bmatrix}$$

where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. This sufficient statistics can be obtained by suitably arranging the terms in the expression

for $p(x_1, \dots, x_n; \boldsymbol{\theta})$:

$$p(x_1, \dots, x_n; \boldsymbol{\theta}) = (2\pi\sigma^2)^{-n/2} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n [(x_i - \bar{x}) + (\bar{x} - \mu)]^2 \right\}$$

and expanding the squares in the exponent.

Example. Suppose x_1, x_2, \dots, x_n , are i.i.d. $\mathcal{N}(\theta, 1)$. Then

$$p(\mathbf{x}; \theta) = \exp\left\{n\theta\left(\bar{x} - \frac{1}{2}\theta\right)\right\} \cdot \underbrace{(2\pi)^{-\frac{1}{2}n} \exp\left\{-\frac{1}{2} \sum_{i=1}^n x_i^2\right\}}_{h(\mathbf{x})}.$$

So, \bar{x} is sufficient by the factorization theorem.

A Side Note on Multivariate Gaussian Pdf

An example (similar to the one in handout # 0). Consider $w[n]$ white Gaussian noise with unit variance:

$$w[n] \sim \underbrace{\mathcal{N}(0, 1)}_{\text{univariate standard normal pdf}}, \quad n = 1, \dots, d$$

implying

$$\begin{aligned} p(w[1], \dots, w[d]) &= \prod_{n=1}^d p(w[n]) \\ &= \frac{1}{(\sqrt{2\pi})^d} \cdot \exp\left(-\frac{1}{2} \cdot \sum_{n=1}^d w[n]^2\right). \end{aligned}$$

This expression can be succinctly written as

$$p(\mathbf{w}) = \underbrace{\frac{1}{(\sqrt{2\pi})^d} \exp\left(-\frac{1}{2} \mathbf{w}^T \mathbf{w}\right)}_{\text{multivariate standard normal pdf}} \triangleq \mathcal{N}(\mathbf{0}, I)$$

where

$$\mathbf{w} = \begin{bmatrix} w[1] \\ w[2] \\ \vdots \\ w[d] \end{bmatrix}.$$

We can generalize:

Definition. A random $d \times 1$ vector \mathbf{X} has a multivariate Gaussian pdf, denoted by

$$\mathbf{X} \sim \mathcal{N}(\underbrace{\boldsymbol{\mu}}_{\mathbb{E}[\mathbf{X}]}, \underbrace{\boldsymbol{\Sigma}}_{\text{cov}(\mathbf{X})})$$

if its pdf is of the form:

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

where $|\boldsymbol{\Sigma}| \equiv$ determinant of $\boldsymbol{\Sigma}$, $\boldsymbol{\mu}$ is a $d \times 1$ vector, and $\boldsymbol{\Sigma}$ is a $d \times d$ symmetric positive definite matrix.

Since $\boldsymbol{\Sigma}$ is symmetric and positive definite, it can be shown that there exists a matrix $\boldsymbol{\Sigma}^{1/2}$ called the square root of $\boldsymbol{\Sigma}$ with the following properties: (i) $\boldsymbol{\Sigma}^{1/2}$ is symmetric, (ii) $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^{1/2} \boldsymbol{\Sigma}^{1/2}$, and (iii) $\boldsymbol{\Sigma}^{1/2} \boldsymbol{\Sigma}^{-1/2} = \mathbf{I}$, where $\boldsymbol{\Sigma}^{-1/2} = (\boldsymbol{\Sigma}^{1/2})^{-1}$.

Theorem. If $\mathbf{W} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and

$$\mathbf{X} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{W}$$

then

$$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

Conversely, if $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then

$$\boldsymbol{\Sigma}^{-1/2} (\mathbf{X} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{0}, I).$$

Suppose that we partition a random $d \times 1$ vector \mathbf{X} as

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_a \\ \mathbf{X}_b \end{bmatrix}.$$

We can partition $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ accordingly:

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}.$$

Theorem. If $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ then

(i) The marginal distribution of \mathbf{X}_a is

$$\mathbf{X}_a \sim \mathcal{N}(\boldsymbol{\mu}_a, \Sigma_{aa}).$$

(ii) The conditional distribution of \mathbf{X}_b given $\mathbf{X}_a = \mathbf{x}_a$ is

$$\mathbf{X}_b | \mathbf{X}_a = \mathbf{x}_a \sim \mathcal{N}\left(\boldsymbol{\mu}_b + \Sigma_{ba} \Sigma_{aa}^{-1} (\mathbf{x}_a - \boldsymbol{\mu}_a), \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}\right).$$

(iii) If \mathbf{a} is a constant $d \times 1$ vector, then

$$\mathbf{a}^T \mathbf{X} \sim \mathcal{N}(\mathbf{a}^T \boldsymbol{\mu}, \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}).$$

(iv)

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \sim \chi_d^2 \quad (\text{Chi-square in your distr. table}).$$

Example: Two Jointly Gaussian RVs

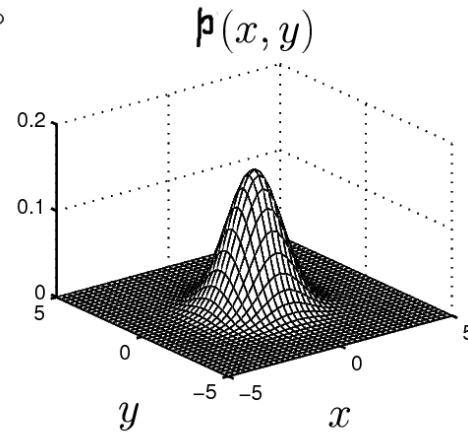
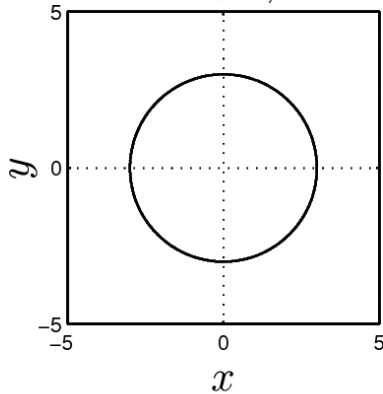
We plot contours of the joint pdf $p_{X,Y}(x,y)$ for zero-mean jointly Gaussian RVs for various values of σ_X, σ_Y , and $\rho_{x,y}$, where we have parametrized $p_{X,Y}(x,y)$ as

$$p_{X,Y}(x,y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho_{X,Y}^2}} \cdot \exp \left\{ -\frac{1}{2(1-\rho_{X,Y}^2)} \cdot \left[\frac{(x-\mu_X)^2}{\sigma_X^2} + \frac{(y-\mu_Y)^2}{\sigma_Y^2} - 2\rho_{X,Y} \frac{(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y} \right] \right\}$$

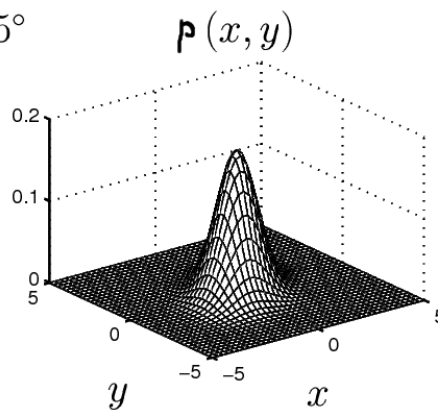
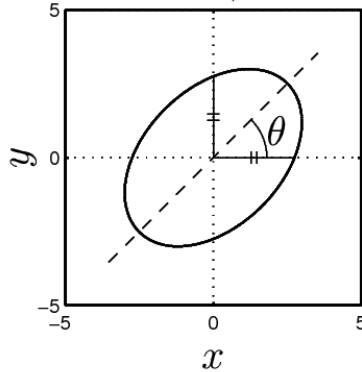
with

$$\underbrace{\rho_{X,Y}}_{\text{correlation coefficient}} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}.$$

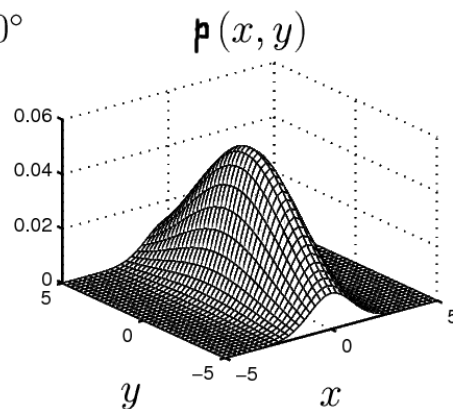
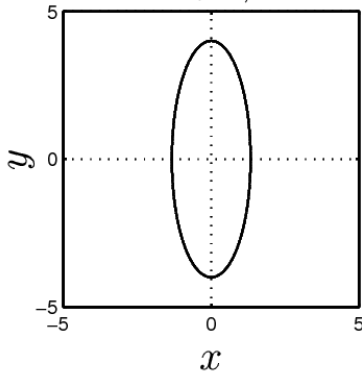
$\sigma_X = 1, \sigma_Y = 1, \rho_{X,Y} = 0: \theta = 0^\circ$



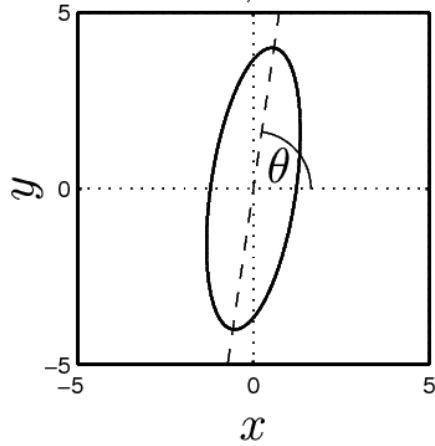
$\sigma_X = 1, \sigma_Y = 1, \rho_{X,Y} = 0.4: \theta = 45^\circ$



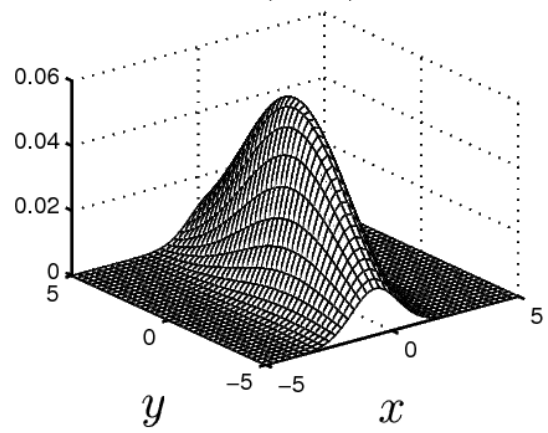
$\sigma_X = 1, \sigma_Y = 3, \rho_{X,Y} = 0: \theta = 90^\circ$



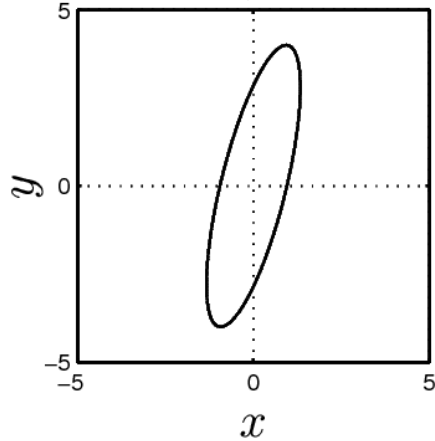
$$\sigma_X = 1, \sigma_Y = 3, \rho_{X,Y} = 0.4: \theta = 81.65^\circ$$



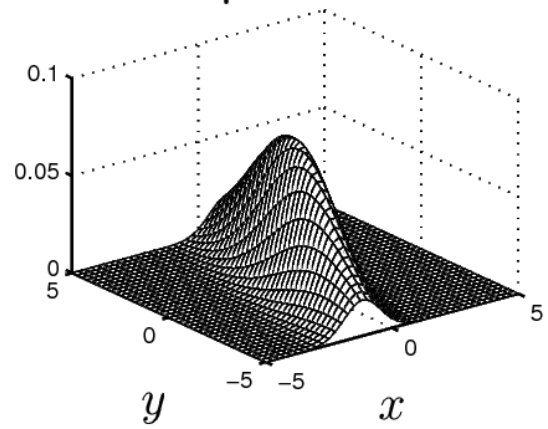
$$p(x, y)$$



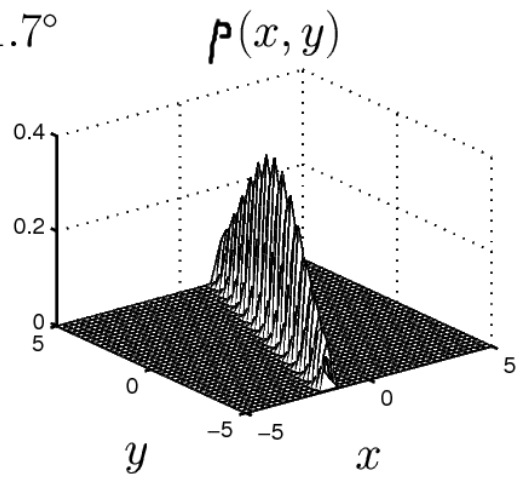
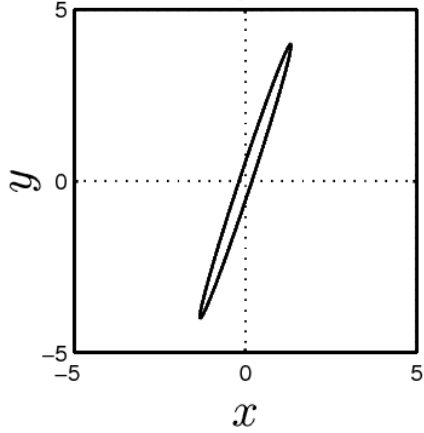
$$\sigma_X = 1, \sigma_Y = 3, \rho_{X,Y} = 0.7: \theta = 76.15^\circ$$



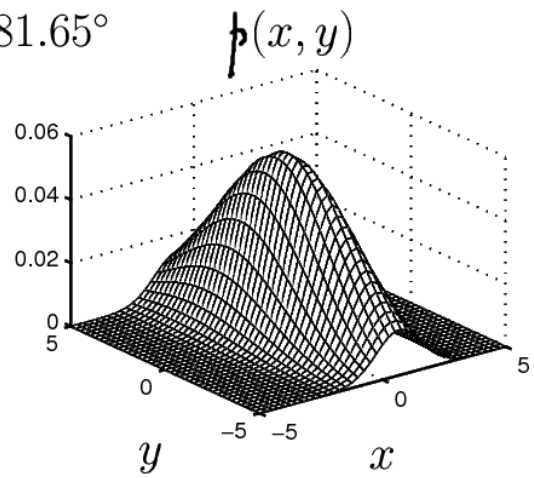
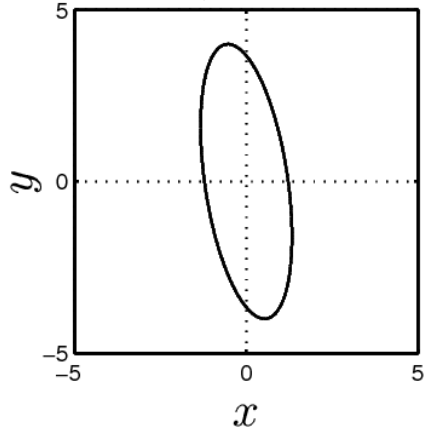
$$p(x, y)$$



$$\sigma_X = 1, \sigma_Y = 3, \rho_{X,Y} = 0.99: \theta = 71.7^\circ$$



$$\sigma_X = 1, \sigma_Y = 3, \rho_{X,Y} = -0.4: \theta = -81.65^\circ$$



Example (Digital Communications)

$$x(t) = \underbrace{s(t)}_{\text{signal}} + \underbrace{w(t)}_{\text{noise}}$$

where the signal $s(t)$ is usually represented using orthonormal basis functions $\varphi_k(t)$:

$$s(t) = \sum_{k=1}^K \alpha_k \varphi_k(t).$$

Note: The signal $s(t)$ is unknown, but it has known structure, incorporated in this basis-function expansion. **We wish to use this structure for data reduction.**

If $\varphi_k(t)$ are orthonormal, it is easy to show that the coefficients α_k can be computed as (for the case of real data):

$$\alpha_k = \int s(t) \varphi_k(t) dt.$$

Here, our goal at the receiver is to decide which $s(t)$ (α_k 's) has been transmitted.

What is typically done in communication receivers is the following: the received data $x(t)$ are matched to the basis

functions, i.e.

$$\hat{\alpha}_k = \int x(t) \varphi_k(t) dt, \quad k = 1, 2, \dots, K \quad (5)$$

are computed and utilized for demodulation.

Question: Are the $\hat{\alpha}_k$ s sufficient statistics for inference about $s(t)$ (or, more precisely, for inference on the α_k s)?

Note: In some applications, sampled data $x[n]$, $n = 0, 1, \dots, N - 1$ are available and

$$\hat{\alpha}_k = \sum_{n=0}^{N-1} x[n] \varphi_k[n], \quad k = 1, 2, \dots, K$$

are used to approximate the integrals in (5) (up to a scaling factor). We focus on this scenario, having in mind that we can easily switch from sums to integrals by letting the sampling interval go to zero — then N will go to infinity. Clearly, N is much larger than the number of basis functions K , i.e.

$$K \ll N.$$

In the sampled-data case, our model is

$$x[n] = \underbrace{s[n]}_{\text{signal}} + \underbrace{w[n]}_{\text{noise}}$$

where

$$s[n] = \sum_{k=1}^K \alpha_k \varphi_k[n].$$

Define

$$\mathbf{x} = \begin{bmatrix} x[0] \\ x[1] \\ \vdots \\ x[N-1] \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w[0] \\ w[1] \\ \vdots \\ w[N-1] \end{bmatrix}, \quad \boldsymbol{\mu}(\boldsymbol{\alpha}) = \begin{bmatrix} s[0] \\ s[1] \\ \vdots \\ s[N-1] \end{bmatrix}$$

and

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_K \end{bmatrix}$$

implying

$$\mathbf{x} = \boldsymbol{\mu}(\boldsymbol{\alpha}) + \mathbf{w}.$$

If the noise is additive zero-mean Gaussian with covariance matrix $\mathbb{E}[\mathbf{w}\mathbf{w}^T] = \mathbf{C}$, then

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}(\boldsymbol{\alpha}), \mathbf{C})$$

which is a multivariate Gaussian pdf:

$$p(\mathbf{x}; \boldsymbol{\alpha}) = \frac{1}{\sqrt{(2\pi)^n |\mathbf{C}|}} \exp \left\{ -\frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\alpha})]^T \mathbf{C}^{-1} [\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\alpha})] \right\} \quad (6)$$

and

$$\boldsymbol{\mu}(\boldsymbol{\alpha}) = \begin{bmatrix} s[0] \\ s[1] \\ \vdots \\ s[N-1] \end{bmatrix} = \begin{bmatrix} \sum_{k=1}^K \alpha_k \varphi_k[0] \\ \sum_{k=1}^K \alpha_k \varphi_k[1] \\ \vdots \\ \sum_{k=1}^K \alpha_k \varphi_k[N-1] \end{bmatrix} = \mathbf{F}\boldsymbol{\alpha}$$

where

$$\mathbf{F} = \begin{bmatrix} \varphi_1(0) & \varphi_2(0) & \cdots & \varphi_K(0) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_1(N-1) & \varphi_2(N-1) & \cdots & \varphi_K(N-1) \end{bmatrix}$$

is an $N \times K$ matrix. So

$$p(\mathbf{x}; \boldsymbol{\alpha}) = \frac{1}{\sqrt{(2\pi)^n |\mathbf{C}|}} \cdot \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mathbf{F}\boldsymbol{\alpha})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{F}\boldsymbol{\alpha}) \right\}.$$

What are the sufficient statistics for inference on $\boldsymbol{\alpha}$? It depends on our knowledge about \mathbf{C} . If \mathbf{C} is *unknown*, we clearly cannot separate out any non-trivial sufficient statistics for both $\boldsymbol{\alpha}$ and \mathbf{C} . If \mathbf{C} is *known*, then the vector of sufficient statistics for $\boldsymbol{\alpha}$ is

$$\mathbf{F}^T \mathbf{C}^{-1} \mathbf{x} \quad (7)$$

which is a $K \times 1$ vector. Since, $K \ll N$, (7) achieves dimensionality reduction compared with the raw data \mathbf{x} . Note

that, if \mathbf{C} is unknown, we cannot compute (7) \implies not realizable.

If $\mathbf{C} = \sigma^2 \mathbf{I}$ (i.e. the noise is white) and the noise variance σ^2 is *known*, then (7) simplifies to (up to a known proportionality factor):

$$\mathbf{F}^T \mathbf{x} = \begin{bmatrix} \sum_{n=0}^{N-1} \varphi_1[n] x[n] \\ \sum_{n=0}^{N-1} \varphi_2[n] x[n] \\ \vdots \\ \sum_{n=0}^{N-1} \varphi_K[n] x[n] \end{bmatrix} = \begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_2 \\ \vdots \\ \hat{\alpha}_K \end{bmatrix}.$$

If σ^2 is *unknown*, then

$$p(\mathbf{x}; \boldsymbol{\alpha}) = \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{x} - \mathbf{F}\boldsymbol{\alpha})^T (\mathbf{x} - \mathbf{F}\boldsymbol{\alpha}) \right\}.$$

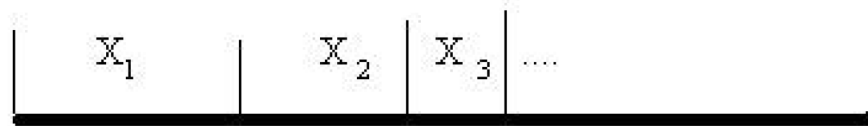
where σ^2 and $\boldsymbol{\alpha} \equiv$ parameters and $\mathbf{x} \equiv$ data. Now

$$\mathbf{x}^T \mathbf{x} = \sum_{n=0}^{N-1} x^2[n] \quad \text{and} \quad \mathbf{F}^T \mathbf{x} = \begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_2 \\ \vdots \\ \hat{\alpha}_K \end{bmatrix}$$

are jointly sufficient for $\boldsymbol{\alpha}$ and σ^2 .

Examples: Computing the Sufficient Statistics

Example: Suppose that elements of $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ are i.i.d. inter-arrival times of packets arriving at a node in a communication network.



We assume that $x_i, i = 1, 2, \dots, n$ come from an exponential $\text{Expon}(\theta)$ distribution, implying

$$p(\mathbf{x}; \theta) = \prod_{i=1}^n \theta \exp(-\theta x_i) = \theta^n \exp\left(-\theta \underbrace{\sum_{i=1}^n x_i}_{T(\mathbf{x})}\right), \quad \forall x_i \geq 0.$$

Example: Elements of $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ are i.i.d. $\text{uniform}(0, \theta)$:

$$p(\mathbf{x}; \theta) = \frac{1}{\theta^n} \prod_{i=1}^n i_{[0, \theta]}(x_i) = \frac{1}{\theta^n} \underbrace{i_{[-\infty, \theta]}(\underbrace{\max x_i}_{T(\mathbf{x})})}_{g(T(\mathbf{x}), \theta)} \cdot \underbrace{i_{[0, \infty]}(\min x_i)}_{h(\mathbf{x})}$$

where $i_A(x)$ denotes the indicator function:

$$i_A(x) = \begin{cases} 1, & x \in A, \\ 0, & \text{otherwise} \end{cases}$$

Here, we have used the facts that

$$x_1, x_2, \dots, x_N < \theta \iff \max x_i < \theta$$

and

$$x_1, x_2, \dots, x_N > 0 \iff \min x_i > 0.$$

Example. Detection problem: $\theta \in \{0, 1\}$ and

$$\begin{aligned} p(x; \theta) &= \theta p(x; 1) + (1 - \theta) p(x; 0) \\ &= \underbrace{\left[\theta \frac{p(x; 1)}{p(x; 0)} + (1 - \theta) \right]}_{g(T(x), \theta)} \underbrace{p(x; 0)}_{h(x)} \end{aligned}$$

$T(x)$ is the likelihood ratio! It is a very useful sufficient statistics because it is one-dimensional regardless of the nature of $p(x; \theta)$. See also Poor, Example IV.C.1.

Definition. The statistic $T(x)$ is *minimally sufficient* if it is sufficient and provides a greater reduction of the data than any other sufficient statistic $S(x)$.

(Multiparameter) Exponential Family of Distributions

$$p(x; \boldsymbol{\theta}) = h(x) \exp \left[\underbrace{\sum_{i=1}^k \eta_i(\boldsymbol{\theta}) T_i(x) - B(\boldsymbol{\theta})}_{\substack{\text{coupling between} \\ \text{parameters and data has} \\ \text{a very specific form}}} \right]$$

By the factorization theorem, $\mathbf{T}(X) = [T_1(X), \dots, T_k(X)]^T$ is sufficient. It is the *natural sufficient statistic* of the family. The exponential family is important:

- It covers quite a few useful distributions, including some that are fairly complex; e.g. Markov random fields, used in image analysis, are virtually all in the exponential-family form;
- It is popular in graphical models as well (Markov random fields again);
- Many methods greatly simplify in the case of exponential family, e.g. the EM algorithm (to be discussed later).

If the support of $p(x; \theta)$ depends on θ , then $p(x; \theta)$ cannot belong to the exponential family. For example, $\text{uniform}(0, \theta)$ is not a member of the exponential family.

Multiple i.i.d. measurements making $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$.
Then,

$$\begin{aligned} p(\mathbf{x}; \theta) &= \prod_{l=1}^n \left\{ h(x_l) \exp \left[\sum_{i=1}^k \eta_i(\boldsymbol{\theta}) T_i(x_l) - B(\boldsymbol{\theta}) \right] \right\} \\ &= \underbrace{\left[\prod_{l=1}^n h(x_l) \right] \cdot \exp \left[\sum_{i=1}^k \eta_i(\boldsymbol{\theta}) \sum_{l=1}^n T_i(x_l) - nB(\boldsymbol{\theta}) \right]}_{\text{again the exponential family}} \end{aligned}$$

and hence the vector of natural sufficient statistics is

$$\mathbf{T}(X) = \left[\sum_{l=1}^n T_1(X_l), \dots, \sum_{l=1}^n T_k(X_l) \right]^T.$$

For more about exponential families, see B & D, Chapter 1.6.

A Side Note on One-Parameter Canonical Exponential Family

Here is a simple special sub-family of the exponential family — *the one-parameter canonical exponential family*:

$$p(\mathbf{x}; \eta) = h(\mathbf{x}) \exp \left[\underbrace{\eta}_{\substack{\text{scalar} \\ \text{canonical} \\ \text{parameter}}} T(\mathbf{x}) - A(\eta) \right]$$

where

$$A(\eta) = \log \int \cdots \int h(\boldsymbol{\chi}) \exp[\eta T(\boldsymbol{\chi})] d\boldsymbol{\chi} \quad \text{for a pdf } p(\mathbf{x}; \eta)$$

$$A(\eta) = \log \sum_{\boldsymbol{\chi}} h(\boldsymbol{\chi}) \exp[\eta T(\boldsymbol{\chi})] \quad \text{for a pmf } p(\mathbf{x}; \eta).$$

If we can compute the normalizing term $A(\eta)$ in a simple form, then it is easy to find the mean and variance of $T(X)$:

$$\mathbb{E}_{p(\mathbf{x}; \eta)}[T(\mathbf{X})] = \frac{dA(\eta)}{d\eta}, \quad \text{var}_{p(\mathbf{x}; \eta)}[T(\mathbf{X})] = \frac{d^2 A(\eta)}{d\eta^2}.$$

Why is this useful? Here is an example. Suppose

x_1, x_2, \dots, x_n are i.i.d. from

$$p(x; \theta) = \underbrace{\frac{x}{\theta^2} \cdot \exp\left(-\frac{x^2}{2\theta^2}\right)}_{\text{Rayleigh pdf}}, \quad x > 0, \theta > 0.$$

Define $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ and write the pdf of \mathbf{x} :

$$\begin{aligned} p(\mathbf{x}; \theta) &= \prod_{i=1}^n \left[\frac{x_i}{\theta^2} \cdot \exp\left(-\frac{x_i^2}{2\theta^2}\right) \right] \\ &= \left(\prod_{i=1}^n x_i \right) \cdot \exp \left[\underbrace{\frac{-1}{2\theta^2}}_{\eta} \cdot \underbrace{\left(\sum_{i=1}^n x_i^2 \right)}_{T(\mathbf{x})} - n \log \theta^2 \right] \end{aligned}$$

implying

$$\theta^2 = -\frac{1}{2\eta}$$

and, consequently,

$$A(\eta) = -n \log(-2\eta).$$

Therefore, the natural sufficient statistic $T(\mathbf{X}) = \sum_{i=1}^n X_i^2$ has mean

$$\mathbb{E} \left[\sum_{i=1}^n X_i^2 \right] = \frac{dA(\eta)}{d\eta} = -\frac{n}{\eta} = 2n\theta^2$$

and variance

$$\text{var} \left[\sum_{i=1}^n X_i^2 \right] = \frac{d^2 A(\eta)}{d\eta^2} = \frac{n}{\eta^2} = 4n \theta^4.$$

Direct computation of these moments would be more complicated.