

Gaussian Derivatives

$$\begin{aligned}\frac{\partial}{\partial x} \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{-x}{\sqrt{2\pi}\sigma^3} e^{-x^2/(2\sigma^2)} = \frac{-x}{\sigma^2} p(x) \\ \frac{\partial^2}{\partial x^2} \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{1}{\sqrt{2\pi}\sigma^5} (-\sigma^2 + x^2) e^{-x^2/(2\sigma^2)} = \frac{-\sigma^2 + x^2}{\sigma^4} p(x) \\ \frac{\partial^3}{\partial x^3} \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{1}{\sqrt{2\pi}\sigma^7} (3x\sigma^2 - x^3) e^{-x^2/(2\sigma^2)} = \frac{-3x\sigma^2 - x^3}{\sigma^6} p(x),\end{aligned}$$

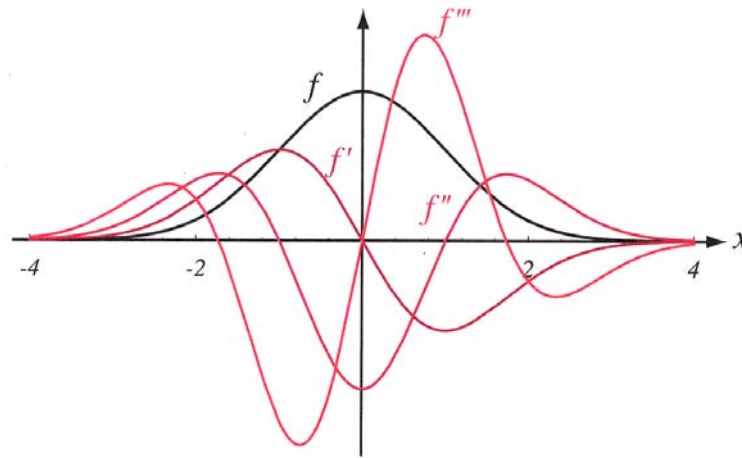


FIGURE A.2. A one-dimensional Gaussian distribution and its first three derivatives, shown for $f(x) \sim N(0, 1)$.

Integral of the Gaussian

$$\text{erf}(\mathbf{u}) = (2 / \sqrt{\pi}) \int_0^{\mathbf{u}} \mathbf{e}^{-\mathbf{x}^2} \mathbf{d}\mathbf{x}$$

Area under the standardized Gaussian

Error Function

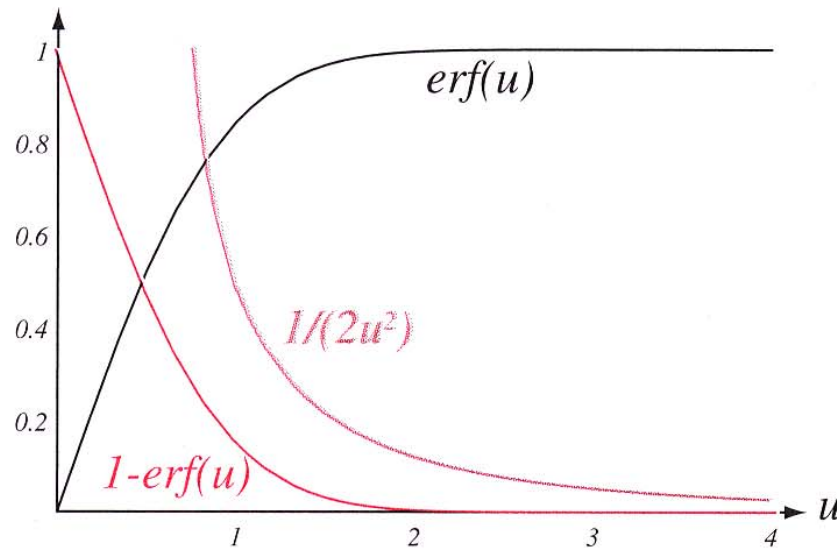


FIGURE A.3. The error function $\text{erf}(u)$ corresponds to the area under a standardized Gaussian between $-\sqrt{2}u$ and $\sqrt{2}u$; that is, if x is a standardized Gaussian random variable, $\Pr[|x| \leq \sqrt{2}u] = \text{erf}(u)$. Thus, the complementary probability, $1 - \text{erf}(u)$, is the probability that a sample is chosen with $|x| > \sqrt{2}u$. Chebyshev's inequality states that for an arbitrary distribution having zero mean and unit standard deviation, $\Pr[|x| > \epsilon]$ is less than $1/\epsilon^2$, so that the lower curve is bounded by $1/(2u^2)$. As shown, this bound is quite loose for a Gaussian.

Generalization to any confidence level

$$error_s(h) \pm z_N \sqrt{\frac{error_s(h)(1 - error_s(h))}{n}}$$

- Values of z_N for 2-Sided $N\%$ Confidence Intervals

Confidence Level $N\%$:	50%	68%	80%	90%	95%	98%	99%
Constant z_n :	0.67	1.00	1.28	1.64	1.96	2.33	2.58

Rule of thumb

- Binomial is approximated by Normal distribution as long as $n \geq 30$ or when $np(1-p) \geq 5$
- For smaller values of n it is wise to use a table giving exact values for the binomial distribution

Bivariate Normal Densities

$$\rho = \frac{\sigma_{12}}{\sigma_1 \sigma_2}.$$

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix},$$

$$|\Sigma| = \sigma_1^2 \sigma_2^2 (1 - \rho^2).$$

$$\begin{aligned} \Sigma^{-1} &= \frac{1}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} \begin{bmatrix} \sigma_2^2 & -\rho\sigma_1\sigma_2 \\ -\rho\sigma_1\sigma_2 & \sigma_1^2 \end{bmatrix} \\ &= \frac{1}{1 - \rho^2} \begin{bmatrix} 1/\sigma_1^2 & -\rho/(\sigma_1\sigma_2) \\ -\rho/(\sigma_1\sigma_2) & 1/\sigma_2^2 \end{bmatrix}. \end{aligned}$$

$$(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$= [(x_1 - \mu_1) (x_2 - \mu_2)] \frac{1}{1 - \rho^2} \begin{bmatrix} 1/\sigma_1^2 & -\rho/(\sigma_1\sigma_2) \\ -\rho/(\sigma_1\sigma_2) & 1/\sigma_2^2 \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}$$

$$= \frac{1}{1 - \rho^2} \left[\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 - 2\rho \left(\frac{x_1 - \mu_1}{\sigma_1} \right) \left(\frac{x_2 - \mu_2}{\sigma_2} \right) + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right].$$

$$p_{x_1 x_2}(x_1, x_2) = \frac{1}{2\pi \sigma_1 \sigma_2 \sqrt{1 - \rho^2}} \quad (112)$$

$$\times \exp \left[-\frac{1}{2(1 - \rho^2)} \left[\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 - 2\rho \left(\frac{x_1 - \mu_1}{\sigma_1} \right) \left(\frac{x_2 - \mu_2}{\sigma_2} \right) + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right] \right].$$

$p(x_1, x_2)$ is a hill-shaped surface

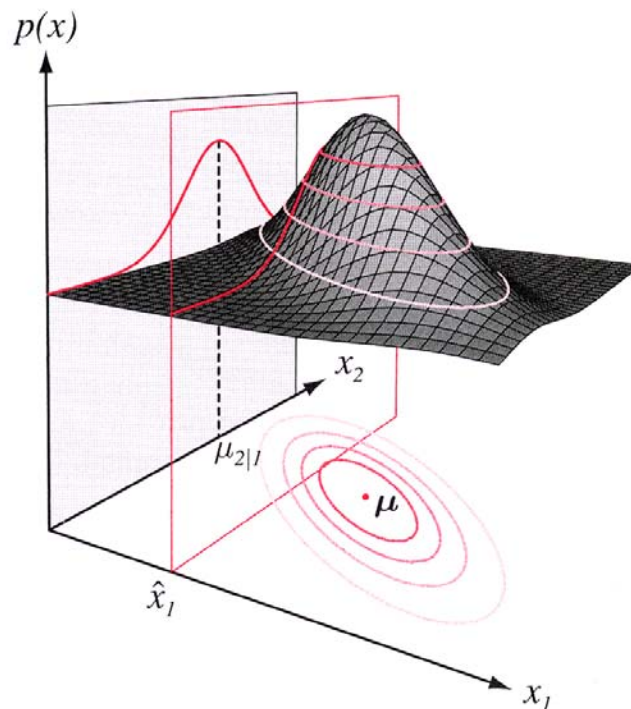


FIGURE A.4. A two-dimensional Gaussian having mean μ and nondiagonal covariance Σ . If the value on one variable is known, for instance $x_1 = \hat{x}_1$, the distribution over the other variable is Gaussian with mean $\mu_{2|1}$.

LEVEL CURVES

- Slice the surface with horizontal planes which the locus of points with the quadratic form

$$\left(\frac{x_1 - \mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x_1 - \mu_1}{\sigma_1}\right)\left(\frac{x_2 - \mu_2}{\sigma_2}\right) + \left(\frac{x_2 - \mu_2}{\sigma_2}\right)^2$$

A.6 HYPOTHESIS TESTING

Statistical hypothesis testing provides a formal way to decide if the results of an experiment are significant or accidental. It is standard statistical terminology to call a set of n measurements $\mathcal{X}_n = \{x_1, x_2, \dots, x_n\}$ a *sample of size n* . However, in keeping with the terminology that is universally used in pattern recognition, we shall call each individual measurement a sample. Suppose that we have a set of samples that are drawn either from a known distribution D_0 or from some other distribution. In pattern classification, we seek to determine which distribution was the source of any sample; and if it is indeed D_0 , we would classify the point accordingly. Hypothesis testing addresses a somewhat different but related problem. We assume initially that distribution D_0 is the source of the patterns; this is called the *null hypothesis* and is often denoted H_0 . Based on the value of any observed sample, we ask whether we can

reject the null hypothesis—that is, state with some degree of confidence (expressed as a probability) that the sample did *not* come from D_0 .

For instance, D_0 might be a standardized Gaussian, $p(x) \sim N(0, 1)$, and hence our null hypothesis is that a sample comes from a Gaussian with mean $\mu = 0$. If the value of a particular sample is small (e.g., $x = 0.3$), it is likely that it came from the D_0 ; after all, 68% of the samples drawn from that distribution have absolute value less than $x = 1.0$ (cf. Fig. A.1). If a sample's value is large (e.g., $x = 5$), then we would be more confident that it did *not* come from D_0 . At such a situation we merely conclude that (with some probability) the sample was drawn from a standardized Gaussian distribution with $\mu \neq 0$.

Viewed another way, for any confidence—expressed as a probability—there exists a criterion value such that if the sampled value differs from $\mu = 0$ by more than that criterion, we reject the null hypothesis. (It is traditional to use confidences of .01 or .05.) We then say that the difference of the sample from 0 is *statistically significant*. For instance, if our null hypothesis is a standardized Gaussian, then if our sample differs from the value $x = 0$ by more than 2.576, we could reject the null hypothesis “at the .01 confidence level,” as can be deduced from Table A.1. A more sophisticated analysis could be applied if *several* samples are all drawn from D_0 or if the null hypothesis involved a distribution other than a Gaussian. Of course, this usage of “significance” applies only to the statistical properties of the problem—it implies nothing about whether the results are “important.” Hypothesis testing is of great generality, and it is useful when we seek to know whether something other than the assumed case (the null hypothesis) is likely to be the case.

A.6.1 Chi-Squared Test

Hypothesis testing can be applied to discrete problems too. Suppose we have n patterns— n_1 of which are known to be in ω_1 , and n_2 in ω_2 —and we are interested in determining whether a particular decision rule is useful or informative. In this case, the null hypothesis is that a random decision rule is present—one that selects a pattern and with some probability P places it in a category which we will call the “left” category, and otherwise in the “right” category. We say that a candidate rule is informative if it differs significantly from such a random decision.

What we need is a clear mathematical definition of statistical significance under these conditions. The random rule (the null hypothesis) would place Pn_1 patterns from ω_1 and Pn_2 from ω_2 independently in the left category and the remainder in the right category. Our candidate decision rule would differ significantly from the random rule if the proportions differed significantly from those given by the random rule. Formally, we let n_{iL} denote the number of patterns from category ω_i placed in the left category by our candidate rule. The so-called *chi-squared* statistic for this case is

$$\chi^2 = \sum_{i=1}^2 \frac{(n_{iL} - n_{ie})^2}{n_{ie}}, \quad (116)$$

where, according to the null hypothesis, the number of patterns in category ω_i that we expect to be placed in the left category is $n_{ie} = Pn_i$. Clearly χ^2 is nonnegative, and it is zero if and only if all the observed numbers n_{iL} match the expected numbers n_{ie} . The higher the value of χ^2 , the less likely it is that the null hypothesis is true. Thus, for a sufficiently high χ^2 , the difference between the expected and observed distributions is statistically significant, we can reject the null hypothesis, and we

Table A.2. Critical Values of Chi-Square (at Two Confidence levels) for Different Degrees of Freedom (df)

df	.05	.01	df	.05	.01	df	.05	.01
1	3.84	6.64	11	19.68	24.72	21	32.67	38.93
2	5.99	9.21	12	21.03	26.22	22	33.92	40.29
3	7.82	11.34	13	22.36	27.69	23	35.17	41.64
4	9.49	13.28	14	23.68	29.14	24	36.42	42.98
5	11.07	15.09	15	25.00	30.58	25	37.65	44.31
6	12.59	16.81	16	26.30	32.00	26	38.88	45.64
7	14.07	18.48	17	27.59	33.41	27	40.11	46.96
8	15.51	20.09	18	28.87	34.80	28	41.34	48.28
9	16.92	21.67	19	30.14	37.57	29	42.56	49.59
10	18.31	23.21	20	31.41	37.57	30	43.77	50.89

can consider our candidate decision rule is “informative.” For any desired level of significance—such as .01 or .05—a table gives the critical values of χ^2 that allow us to reject the null hypothesis (Table A.2).

There is one detail that must be addressed: the number of degrees of freedom (df). In the situation described above, once the probability P is known, there is only one free variable needed to describe a candidate rule. For instance, once the number of patterns from ω_1 placed in the left category are known, all other values are determined uniquely. Hence in this case the number of degrees of freedom is 1. If there were more categories, or if the candidate decision rule had more possible outcomes, then the number of degrees of freedom df would be greater than 1. The higher the number of degrees of freedom, the higher must be the computed χ^2 to meet a desired level of significance.

We denote the critical values as, for instance, $\chi^2_{.01(1)} = 6.64$, where the subscript denotes the significance, here .01, and the integer in parentheses is the number of degrees of freedom. (In Table A.2, we conform to the usage in statistics, where this positive integer is denoted df , despite the possible confusion in calculus where it denotes an infinitesimal real number.) Thus if we have one degree of freedom and the observed χ^2 is greater than 6.64, then we can reject the null hypothesis and say that at the .01 confidence level our results did not come from a (weighted) random decision.

INFORMATION THEORY

$$H = - \sum_{i=1}^m P_i \log_2 P_i,$$

$$H = - \int_{-\infty}^{\infty} p(x) \ln p(x) dx,$$

$$\delta(x - a) = \begin{cases} 0 & \text{if } x \neq a \\ \infty & \text{if } x = a, \end{cases} \quad \text{with}$$

$$\int_{-\infty}^{\infty} \delta(x) dx = 1,$$

Relative Entropy

$$D_{KL}(p(x), q(x)) = \sum_x q(x) \ln \frac{q(x)}{p(x)}.$$

$$D_{KL}(p(x), q(x)) = \int_{-\infty}^{\infty} q(x) \ln \frac{q(x)}{p(x)} dx.$$

Mutual Information

$$I(p; q) = H(p) - H(p|q) = \sum_{x,y} r(x, y) \log_2 \frac{r(x, y)}{p(x)q(y)},$$

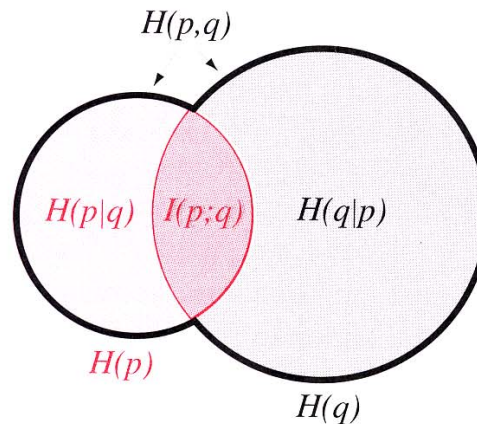


FIGURE A.5. For two distributions p and q , this figure shows the mathematical relationships among the entropy, mutual information $I(p; q)$, and conditional entropies $H(p|q)$ and $H(q|p)$. For instance $I(p; p) = H(p)$; if $I(p; q) = 0$, then $H(q|p) = H(q)$; $H(p, q) = H(p|q) + H(q)$; and so forth.