

Introduction to Neural Networks

U. Minn. Psy 5038

More probability

■ Initialize standard library files:

```
Off[General::spell1];  
SetOptions[ContourPlot, ImageSize -> Small];  
SetOptions[Plot, ImageSize -> Small];  
SetOptions[ListPlot, ImageSize -> Small];
```

Goals

Review the basics of probability distributions and statistics

More on generative modeling: drawing samples

Graphical models for inference

Optimal inference and Task dependence

Probability overview

Random variables, discrete probabilities, probability densities, cumulative distributions

■ Discrete: random variable X can take on a finite set of discrete values

$X = \{x(1), \dots, x(N)\}$

$$\sum_{i=1}^N p_i = \sum_{i=1}^N p(X = x(i)) = 1$$

■ Densities: X takes on continuous values, x , in some range.

Density : $p(x)$

Analogous to material mass,

we can think of the probability over some small domain of the random variable as "probability mass" :

$$\text{prob}(x < X < dx + x) = \int_x^{x+dx} p(x) dx$$

$$\text{prob}(x < X < dx + x) \approx p(x) dx$$

With the mass analogy, however, an object (event space) always "weighs" 1 :

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

Cumulative distribution:

$$\text{prob}(X < x) = \int_{-\infty}^x p(X) dX$$

■ Densities of discrete random variables

The Dirac Delta function, $\delta[\bullet]$, allows us to use the mathematics of continuous distributions for discrete ones, by defining the density as:

$$p[x] = \sum_{i=1}^N p_i \delta[x - x[i]], \text{ where } \delta[x - x[i]] = \begin{cases} \infty & \text{for } x = x[i] \\ 0 & \text{for } x \neq x[i] \end{cases}$$

Think of the delta function, $\delta[\bullet]$, as ϵ wide and $1/\epsilon$ tall, and then let $\epsilon \rightarrow 0$, so that:

$$\int_{-\infty}^{\infty} \delta(y) dy = 1$$

The density, $p[x]$, is a series of spikes. It is infinitely high only at those points for which $x = x[i]$, and zero elsewhere. But "infinity" is scaled so that the local mass or area around each point $x[i]$, is p_i .

■ Joint probabilities

Prob (X AND Y) = $p(X, Y)$

Joint density : $p(x, y)$

Three basic rules of probability

Suppose we know everything there is to know about a set of variables (A,B,C,D,E). What does this mean in terms of probability? It means that we know the joint distribution, $p(A,B,C,D,E)$. In other words, for any particular combination of values (A=a,B=b, C=c, D=d,E=e), we can calculate, look up in a table, or determine some way or another the number $p(A=a,B=b, C=c, D=d,E=e)$, for any particular instances, a, b, c, d, e.

■ Rule 1: Conditional probabilities from joints: The product rule

Probability about an event changes when new information is gained.

Prob(X given Y) = $p(X|Y)$

$$p(X | Y) = \frac{p(X, Y)}{p(Y)}$$

$$p(X, Y) = p(X | Y) p(Y)$$

The form of the product rule is the same for densities as for probabilities.

■ Rule 2: Lower dimensional probabilities from joints: The sum rule (marginalization)

$$p(X) = \sum_{i=1}^N p(X, Y(i))$$

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

■ Rule 3: Bayes' rule

From the product rule, and since $p[X,Y] = p[Y,X]$, we have:

$$p(Y | X) = \frac{p(X | Y) p(Y)}{p(X)}, \text{ and using the sum rule, } p(Y | X) = \frac{p(X | Y) p(Y)}{\sum_Y p(X, Y)}$$

■ Bayes Terminology in inference

Suppose we have some partial data (see half of someone's face), and we want to recall or complete the whole. Or suppose that we hear a voice, and from that visualize the face. These are both problems of statistical inference. We've already studied how to complete a partial pattern using energy minimization, and how energy minimization can be viewed as probability maximization.

We typically think of the **Y** term as a random variable over the hypothesis space (a face), and **X** as data or a stimulus (partial face, or sound). So for recalling a pattern **Y** from an input stimulus **X**, We'd like to have a function that tells us:

$p(\mathbf{Y} | \mathbf{X})$ which is called the **posterior** probability of the hypothesis (e.g. description of the full face as output) given the stimulus (partial face as "data").

-- i.e. what you get when you condition the joint by the probability of the stimulus data. The posterior is often what we'd like to base our decisions on, because it can be proved that picking the hypothesis **Y** which maximizes the posterior (i.e. maximum a posteriori or **MAP** estimation) minimizes the average probability of error.

$p(\mathbf{Y})$ is the **prior** probability of the hypothesis. Some hypotheses are "a priori" more likely than others. But even if it isn't made explicit, a model prior implicitly assumes conditions. Given a context, such as your room, some faces are more likely than others. For me an image patch stimulating my retina in my kitchen is much more likely to be my wife's than my brother's (who lives in another state). Priors are contingent, i.e. conditional on context, $p(\mathbf{Y} | \text{context})$, even if the context is not made explicitly.

$p(\mathbf{X} | \mathbf{Y})$ is the **likelihood** of the hypothesis. Note this is a probability of **X**, but not of **Y**. (The sum over X is one, but the sum over Y isn't necessarily one.)

■ Independence

Knowledge of one event doesn't change the probability of another event.

$p(\mathbf{X}) = p(\mathbf{X} | \mathbf{Y})$ which by the product rule is:

$$p(\mathbf{X}, \mathbf{Y}) = p(\mathbf{X})p(\mathbf{Y})$$

Deterministic relationships

Deterministic relationships can be treated as special cases, and provide a useful way to build some intuitions about probabilities.

For example, suppose we know that $Y = X^2$ exactly, for integer values of $0 < X < 5$. What is the probability of $X = x$, $Y = y$, over the space of possible x 's?

```
In[38]:= p[y_, x_] := If[0 < x < 5, KroneckerDelta[y - x^2] / 4, Null]
```

```
In[39]:= Table[p[y, x], {x, 1, 4}, {y, 1, 16, 1}] // MatrixForm
```

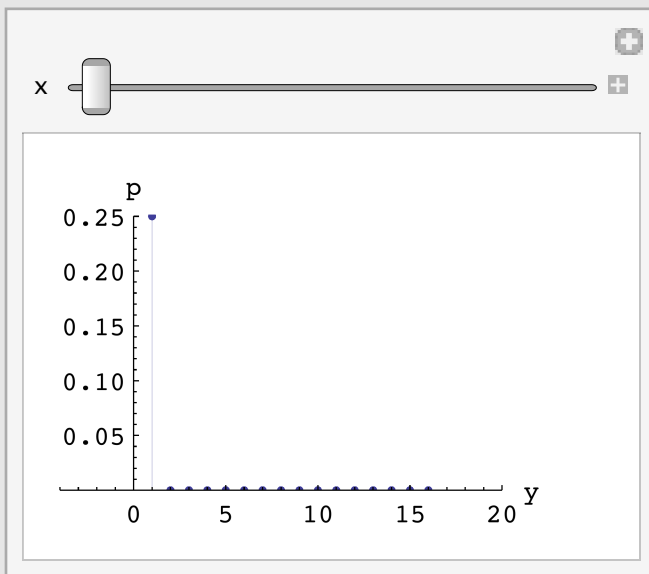
```
Out[39]//MatrixForm=
```

$$\begin{pmatrix} \frac{1}{4} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{4} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{4} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{4} \end{pmatrix}$$

Given $X=x$, what is y ?

```
In[193]:= Manipulate[ListPlot[Table[p[y, x], {y, 1, 16}], Filling -> Axis,
  ImageSize -> Small, Axes -> {True, True}, PlotRange -> {{-4, 20}, {0, .25}},
  AxesLabel -> {"y", "p"}], {x, 1, 16, 1}]
```

```
Out[193]=
```



Note that we've plotted $p(y, X=2)$. What is $p(y | X = 2)$?

What is $p(y)$?

```
In[198]:= py[y_] := Sum[p[y, x], {x, 1, 4, 1}]
```

```
In[202]:= py[4]
```

```
Out[202]=
```

$$\frac{1}{4}$$

```
In[208]:= p[9, 2]
```

```
Out[208]= 0
```

Density mapping theorem

Suppose we have a change of variables that maps a discrete set of x 's uniquely to y 's: $X \rightarrow Y$.

■ Discrete random variables

No change to probability function. The mapping just corresponds to a change of labels, so the probabilities $p(X)=p(Y)$.

■ Continuous random variables

Form of probability density function does change because we require the probability "mass" to be unchanged: $p(x)dx = p(y)dy$

Suppose, $y=f(x)$

$$p_Y(Y) \delta Y = p_X(X) \delta X$$

(In higher dimensions, the transformation is done by multiplying the density by the Jacobian, the determinant of the matrix of partial derivatives of the change of coordinates.)

One can express the density mapping theorem as:

$$p_Y(y) = \int \delta(y - f(x)) f^{-1}(x) p_X(x) dx$$

over each monotonic part of f .

Transformation of variables is used in making random number generators for probability densities other than the uniform distribution, such as a Gaussian.

Convolution theorem for adding rvs

Let x be distributed as $g(x)$, and y as $h(x)$. Then the probability density for $z=x+y$ is, $f(z)$:

$$f(z) = \int g(s) h(z-s) ds \quad (1)$$

Statistics

■ Expectation & variance

Analogous to center of mass:

Definition of expectation or average:

$$\text{Average}[X] = \bar{X} = E[X] = \sum x[i] p[x[i]] \sim \sum_{i=1}^N x_i / N$$

$$\mu = E[X] = \int x p(x) dx$$

Some rules:

$$E[X+Y] = E[X] + E[Y]$$

$$E[aX] = aE[X]$$

$$E[X+a] = a + E[X]$$

Definition of variance:

$$\sigma^2 = \text{Var}[X] = E[(X-\mu)^2] = \sum_{j=1}^N (p(x(j))) (x(j) - \mu)^2 = \sum_{j=1}^N p_j (x_j - \mu)^2$$

$$\text{Var}[X] = \int (x - \mu)^2 p(x) dx \sim \sum_{i=1}^N (x_i - \mu)^2 / N$$

Standard deviation:

$$\sigma = \sqrt{\text{Var}[X]}$$

Some rules:

$$\text{Var}[X] = E[X^2] - E[X]^2$$

$$\text{Var}[aX] = a^2 \text{Var}[X]$$

■ Covariance & Correlation

Covariance:

$$\text{Cov}[X, Y] = E[(X - \mu_X) (Y - \mu_Y)]$$

Correlation coefficient:

$$\rho[X, Y] = \frac{\text{Cov}[X, Y]}{\sigma_X \sigma_Y}$$

■ Covariance matrix

Suppose now that X is a vector: $\{X_1, X_2, \dots\}$

Then we can describe the covariance between pairs of elements of X :

$$\Sigma_{ij} = \text{cov}[X_i, X_j] = E[(X_i - \mu_{X_i})(X_j - \mu_{X_j})] \sim \frac{\sum_{n=1}^N (x_i^n - \mu_{X_i})(x_j^n - \mu_{X_j})}{N}$$

In matrix form, the covariance can be written:

$$\Sigma = \text{cov}[X] = E[(X - E[X])(X - E[X])^T]$$

In other words, the covariance matrix can be approximated by the average outer product. In the language of neural networks, it is a Hebbian matrix memory of pair-wise relationships.

■ Independent random variables

If $p(X, Y) = p(X)p(Y)$, then

$$E[XY] = E[X]E[Y] \quad (\text{uncorrelated})$$

$$\text{Cov}[X, Y] = \rho[X, Y] = 0$$

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$$

If two random variables are uncorrelated, they are not necessarily independent.

Two random variables are said to be orthogonal if their correlation is zero.

Degree of belief vs., relative frequency

What is the probability that the Vikings will win the Superbowl in 2010? The Packers? Assigning a number between 0 and 1 is assigning a degree of belief. These probabilities are also called subjective probabilities.

What is the probability that a coin will come up heads? In this case, we can do an experiment. Flip the coin n times, and count the number of heads, say $h[n]$, and then set the probability, $p = h[n]/n$ -- the relative frequency. Of course, if we did it again, we may not get the same estimate of p . One solution often given is:

$$p = \lim_{n \rightarrow \infty} \frac{h(n)}{n}$$

A problem with this, is that in general there is no guarantee that a well – defined limit exists.

In some domains we can measure statistics, and model probabilities of both inputs and outputs. So the relative frequency interpretation seems reasonable. In practice, the dimensions of many problems in perception, cognition, language, and memory are so high, that it is impractical to do this. Suppose you wanted to estimate the joint probabilities of all 6 letter english words (or worse yet, 8x8 pixel images). There are over 300 million possible combinations of 26 letters--i.e. over 300 million "bins" for your word counts. Most of these would have zero or near-zero entries, and it would be hard to get good estimates of most of the joint probabilities of 6 letter combinations. Although there are ways to estimate "objective priors" in high-dimensional spaces (see below), once we use the statistical framework to model perception, say of a particular cue (say), then probabilities can become more like "subjective unconscious beliefs". From a modeling perspective, one can treat a specific prior as an assumed ingredient, and test to see how well the model accounts for the data, and

how well it predicts new data. If the model does a poor job empirically, then one can go back and question whether an alternative prior could improve the predictions.

Principle of insufficient reason

■ Principle of symmetry

Suppose we have N events, $x[1], x[2], x[3], \dots, x[N]$ that are all physically identical except for the label. Then assume that

$$\text{prob}(x(1)) = \text{prob}(x(2)) = \text{prob}(x(3)) = \text{prob}(x(N)) = \frac{1}{N}$$

In other words, if we have no additional information about the events, we should assume that they are uniformly distributed. I.e., assume a *uniform prior*.

What about the continuous case where there is no reason to assume any particular value at all between $-\infty$ and $+\infty$?

Improper priors.

■ Information theory and Maximum entropy

Information theory provides a powerful extension to the principle of symmetry. Information of event (or signal) X is:

$$\text{Information}[X] = -\log_2(p(X))$$

The more improbable an event, the more surprising it is in some sense, and the more information it provides.

Using the definition of expectation above, we can specify the expectation of information (or average information), which is called entropy. Entropy of a random variable X with probability distribution $p[X]$ is:

$$H(X) = \text{Average}(\text{Information}[X]) = -\sum_X p(X) \log_2(p(X))$$

It can be shown that out of all possible probability distributions, $H(X)$ is biggest for the uniform distribution, $p(X)=1/N$. Maximum entropy is looking like the symmetry principle.

Indeed it turns out that a more powerful formulation of the principle of symmetry is maximum entropy. For example, out of all possible probability distributions of a random variable with infinity range, but with a specific mean and standard deviation, the Gaussian is unique in having the largest entropy. If the range goes from zero to infinity, and we know the mean, the maximum entropy distribution is an exponential (Cover and Thomas).

An interesting application of the maximum entropy principle is to learning image textures joint probabilities: $p(I[1], \dots, I[N])$, where N is very big, but where one has only a relatively small number of measured statistics relative to the number of possible images (which is really huge). The measurements underdetermine the dimensionality of the probability space--i.e. there are many different probability distributions that give the same statistics. So the principle of symmetry, or insufficient reason, says to choose the one with the maximum entropy. This provides a way to model high dimensional priors. And in fact, this has been done for texture models briefly described in a previous lecture (Zhu et al.).

Examples of sampling, univariate distributions, multivariate gaussian, mixtures

Making a univariate (scalar) gaussian random number generator:

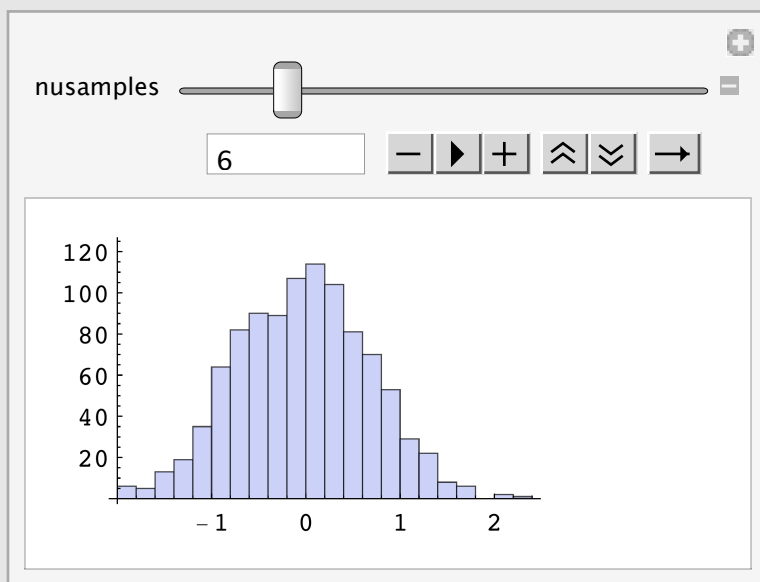
■ Method 1: Just for Gaussian. Use Central Limit Theorem

If all we want to do is make a Gaussian random number generator from a uniformly distributed generator, we can use the Central Limit Theorem.

Try the cell below with `nusamples = 1, 2, ..., 10, ...`

```
In[125]:= Manipulate[
  z1 = Table[
    
$$\sum_{i=1}^{\text{nusamples}} \text{RandomReal}[] - \frac{\text{nusamples}}{2}$$
,
    {1000}];
  Histogram[z1, ImageSize -> Small], {nusamples, 1, 30, 1}]
```

Out[125]=



Method 2: Use Density Mapping theorem. More general.

We'll use the density mapping theorem to turn uniformly distributed random numbers `RandomReal[]` into gaussian distributed random numbers with mean =0 and standard deviation =1.

$$p_Y (Y) \delta Y = p_X (X) \delta X$$

$$p_Y (Y) \frac{\delta Y}{\delta X} = p_X (X)$$

Suppose $p_Y (Y) = 1$ (over the unit interval, but zero elsewhere). Then

$$Y (X) = \int_{-\infty}^X p_X (X') dX' = P (X) \quad (2)$$

Thus if we sample from the uniform distribution to get y , x should be distributed according to $p_X (X)$. To do this, we need a mapping from $y \rightarrow x$. This is given by the inverse cumulative distribution, i.e. $P^{-1}(y)$.

Let's implement this. The quick way is to use *Mathematica's* built-in function to get the inverse cumulative.

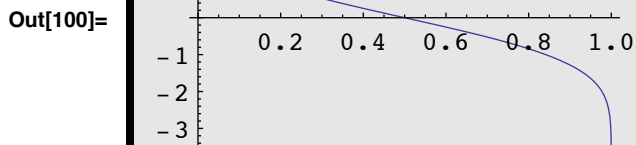
Method 2a: Applied to Gaussian

`InverseErf[]` is the inverse of :

$$\text{erf} (z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$$

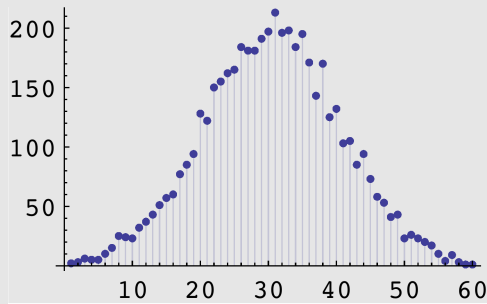
We can use this to define a function for the inverse cumulative of a gaussian:

```
In[98]:= Clear[z];
z[p_] := Sqrt[2] InverseErf[1 - 2 p];
Plot[z[y], {y, 0, 1}]
```



```
In[101]:= binsize = 0.1;
z1 = Table[z[RandomReal[]], {5000}];
freq = BinCounts[z1, {-3, 3, binsize}];
ListPlot[freq, Filling -> Axis]
```

Out[104]=



Method 2b: From scratch: Works for almost any distribution.

Suppose we have a discrete representation of any cumulative distribution. How can we generate samples? For illustration purposes, we'll illustrate the method with a discretization of the Gaussian.

Our first goal is to produce a discrete approximation to the cumulative gaussian. To review where things come from, we'll start with the definition of a Gaussian, and make sure it is normalized.

```
In[126]:= Integrate[Exp[-(x - x0)^2 / (2 * σ^2)], {x, -Infinity, Infinity}]
```

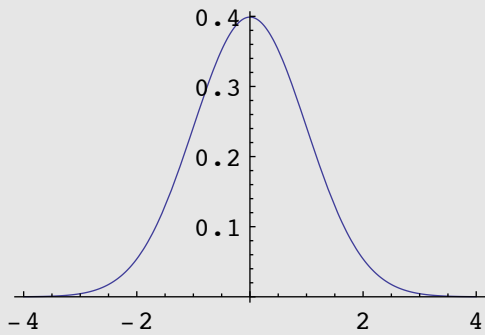
Out[126]=

$$\text{If}\left[\text{Re}\left[\sigma^2\right] > 0, \frac{\sqrt{2\pi}}{\sqrt{\frac{1}{\sigma^2}}}, \text{Integrate}\left[e^{-\frac{(x-x_0)^2}{2\sigma^2}}, \{x, -\infty, \infty\}, \text{Assumptions} \rightarrow \text{Re}\left[\sigma^2\right] \leq 0\right]\right]$$

Let $x_0=0$ and $\sigma=1$:

```
In[127]:= Plot[ $\frac{e^{-\frac{x_1^2}{2}}}{\sqrt{2\pi}}$ , {x1, -4, 4}]
```

```
Out[127]=
```



Note that `Plot[PDF[NormalDistribution[0,1],x1],{x1,-4,4}]`; gives the same thing using the add-on normal distribution function.

■ Cumulative gaussian

```
In[128]:= Clear[cumulgauss, x, x1];
cumulgauss[x_] := NIntegrate[Exp[-(x1^2)/2]/(Sqrt[2*Pi]),
{x1, -Infinity, x}]
```

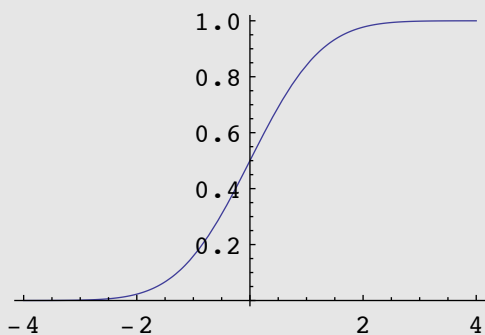
```
In[130]:= cumulgauss[Infinity]
```

```
Out[130]= 1.
```

We can plot up cumulgauss (not very efficiently):

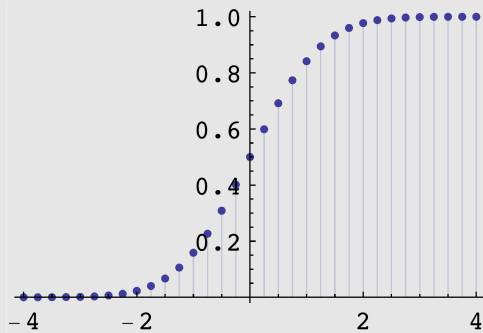
```
In[131]:= Plot[cumulgauss[x], {x, -4, 4}]
```

```
Out[131]=
```



```
In[132]:= lcumulgauss = Table[{x, cumulgauss[x]}, {x, -4., 4., 0.25.}];
ListPlot[lcumulgauss, Filling -> Axis]
```

Out[133]=



■ Make inverse cumulative gaussian table

This is a useful trick whenever you want an inverse function, given a discrete representation.

```
In[134]:= invlcumulgauss = RotateLeft[lcumulgauss, {0, 1}];
```

To see what this does, evaluate:

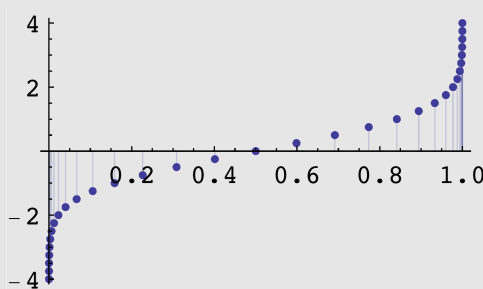
```
In[135]:= {{x1, y1}, {x2, y2}, {x3, y3}}
RotateLeft[{{x1, y1}, {x2, y2}, {x3, y3}}, {0, 1}]
```

Out[135]= {{x1, y1}, {x2, y2}, {x3, y3}}

Out[136]= {{y1, x1}, {y2, x2}, {y3, x3}}

```
In[137]:= ListPlot[invlcumulgauss, Filling -> Axis]
```

Out[137]=



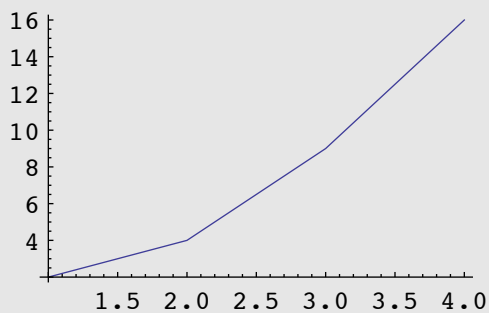
■ Make interpolated function of the inverse cumulative

Another useful trick.

Interpolation works by fitting polynomial curves to the data. Try the test below with various interpolation orders (the default is 3)

```
In[138]:= test = Interpolation[{{1, 2.}, {2, 4}, {3, 9}, {4, 16.}},  
    InterpolationOrder -> 1];  
Plot[test[x], {x, 1, 4}]
```

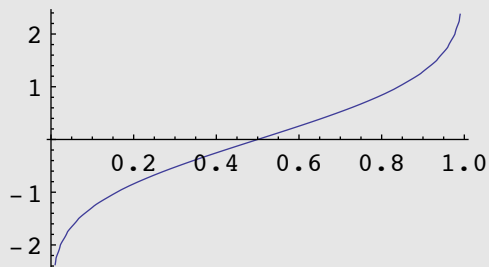
Out[139]=



```
In[140]:= interinvcumulgauss = Interpolation[invcumulgauss];
```

```
In[141]:= Plot[interinvcumulgauss[x], {x, 0.01, 0.99}]
```

Out[141]=



■ Draw samples with a standard deviation of Sqrt[10]

```
In[211]:= Round[10 interinvcumulgauss[RandomReal[]]]
```

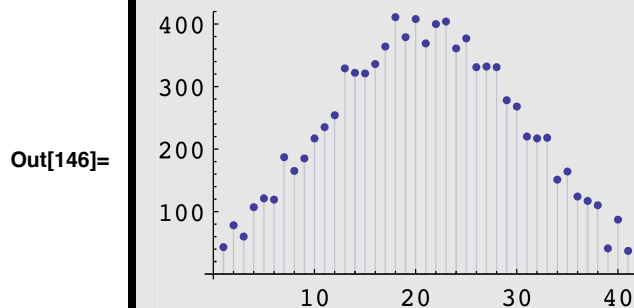
Out[211]=

12

Draw a bunch of samples, and plot up histogram

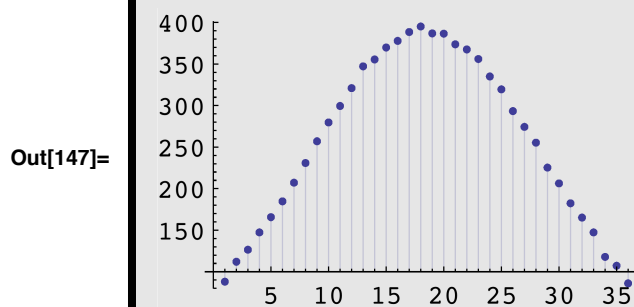
```
In[143]:= z = Table[Round[10 Interinvlcumulgauss[RandomReal[]]], {10 000}];
domain = Range[-20, 20];
Freq = (Count[z, #1] &) /@ domain;
```

```
In[146]:= ListPlot[Freq, Filling -> Axis]
```



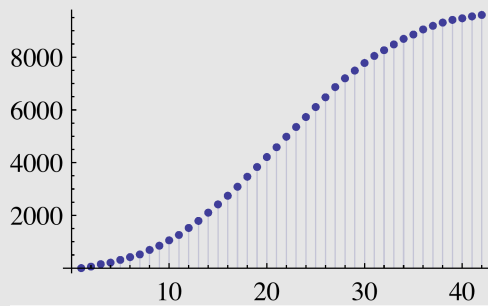
digression...a quick & dirty way to smooth is to do a moving average

```
In[147]:= ListPlot[MovingAverage[Freq, 6], Filling -> Axis]
```



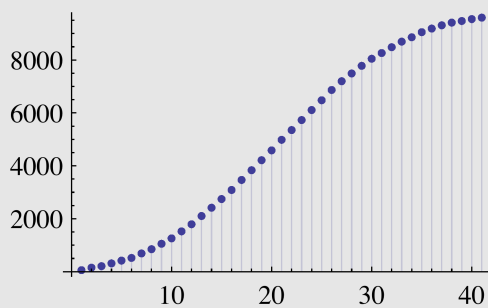
■ Plot up cumulative histogram

```
CumFreq = FoldList[Plus, 0, Freq];
ListPlot[CumFreq, Filling -> Axis]
```



Same thing, with Accumulate[] (new in *Mathematica* 6):

```
CumFreq = Accumulate[Freq];
ListPlot[CumFreq, Filling -> Axis]
```



Multivariate (vector) gaussian distributions

■ Define multivariate gaussian probability density

An n -variate multivariate gaussian (multinormal) distribution with mean vector μ and covariance matrix Σ is denoted $N_n(\mu, \Sigma)$. The density is:

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} \text{Det}[\Sigma]^{1/2}} \text{Exp}\left[-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right] \quad (3)$$

You can read in the *Mathematica* package which has predefined **MultinormalDistribution** $[\mu, \Sigma]$.

```
In[212]:= Needs["MultivariateStatistics`"]
```

Then, for example, you can define the probability density function for mean vector $\{\mu_1, \mu_2\}$, and covariance matrix $\{\{\sigma_{11}^2, \rho * \sigma_{11} * \sigma_{22}\}, \{\rho * \sigma_{11} * \sigma_{22}, \sigma_{22}^2\}\}$, where ρ parameterizes correlation.

```
In[214]:= PDF[MultinormalDistribution[{μ1, μ2},
  {{σ11^2, ρ * σ11 * σ22}, {ρ * σ11 * σ22, σ22^2}}], {x, y}]
```

```
Out[214]= 
$$\frac{e^{\frac{1}{2} \left( -(y-\mu_2) \left( \frac{(y-\mu_2) \sigma_{11}^2}{\sigma_{11}^2 \sigma_{22}^2 - \rho^2 \sigma_{11}^2 \sigma_{22}^2} - \frac{\rho (x-\mu_1) \sigma_{11} \sigma_{22}}{\sigma_{11}^2 \sigma_{22}^2 - \rho^2 \sigma_{11}^2 \sigma_{22}^2} \right) - (x-\mu_1) \left( -\frac{\rho (y-\mu_2) \sigma_{11} \sigma_{22}}{\sigma_{11}^2 \sigma_{22}^2 - \rho^2 \sigma_{11}^2 \sigma_{22}^2} + \frac{(x-\mu_1) \sigma_{22}^2}{\sigma_{11}^2 \sigma_{22}^2 - \rho^2 \sigma_{11}^2 \sigma_{22}^2} \right) \right)}}{2 \pi \sqrt{\sigma_{11}^2 \sigma_{22}^2 - \rho^2 \sigma_{11}^2 \sigma_{22}^2}}$$

```

...but let's do it from scratch. We define a 2-variate density:

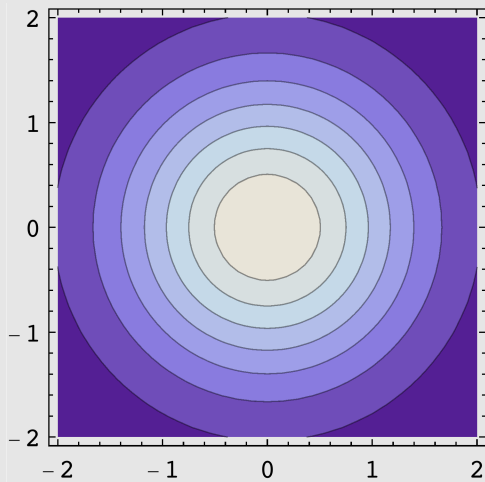
```
In[148]:= multigauss[x_, m_, cov_] :=
Module[{IC, detCov, norm, p},
  IC = Inverse[cov];
  detCov = Abs[Det[cov]];
  norm = N[Sqrt[(2Pi)^2 detCov]];
  p = Exp[-0.5 (x-m).IC.(x-m)]/norm;
  Return[p];
];
```

■ Two variable examples

■ Zero mean, zero correlation

```
In[149]:= m1 = {0, 0}; Cov = {{1, 0}, {0, 1}};  
ContourPlot[multigauss[{x1, x2}, m1, Cov], {x1, -2, 2}, {x2, -2, 2}]
```

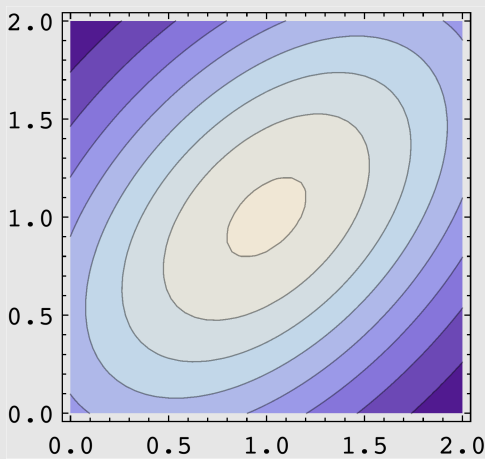
Out[149]=



■ Mean = {1,1}, positive correlation

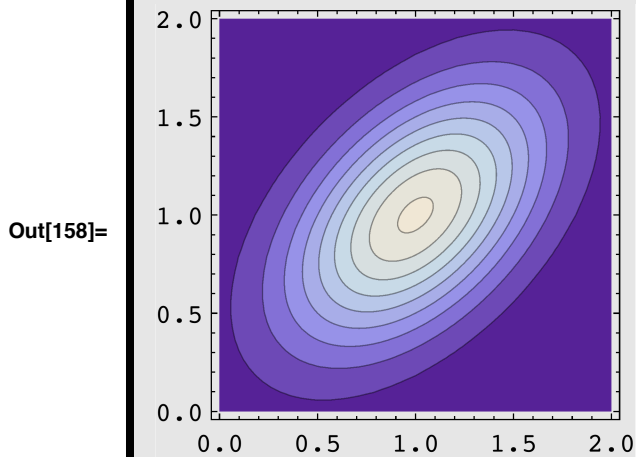
```
In[157]:= m1 = {1, 1}; Cov = {{1, 0.5`}, {0.5`, 1}};  
ContourPlot[multigauss[{x1, x2}, m1, Cov], {x1, 0, 2}, {x2, 0, 2}]
```

Out[157]=



■ Mean = {1,1}, positive correlation, small variance

```
In[158]:= m1 = {1, 1}; Cov = 0.2 {{1, 0.5}, {0.5, 1}};
ContourPlot[multigauss[{x1, x2}, m1, Cov], {x1, 0, 2}, {x2, 0, 2}]
```



■ Mixture of gaussians

α is a mixing parameter

$$p(x) = \alpha p_1(x) + (1 - \alpha) p_2(x) \text{ where } 0 \leq \alpha \leq 1 \quad (4)$$

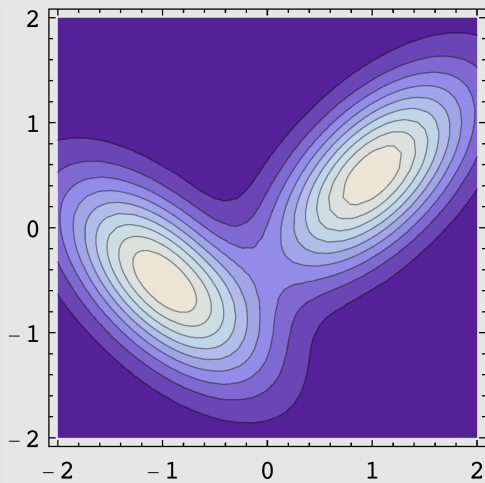
α can be interpreted in terms of a prior probability of choosing which of two distributions a sample will be drawn from.

Starting with $p(x,a)$ use the sum and product rules to derive the above mixture density where $p(a=a_1) = \alpha$, and $p(a = a_2) = (1-\alpha)$.

```
In[159]:= m1 = {1, .5}; m2 = {-1, -.5};
Cov1 = 0.4*{{1, .6}, {.6, 1}};
Cov2 = 0.4*{{1, -.6}, {-.6, 1}};
mix[x_] := 0.5 (multigauss[x,m1,Cov1] + multigauss[x,m2,Cov2]);
```

```
In[163]:= ContourPlot[mix[{x1, x2}], {x1, -2, 2}, {x2, -2, 2}]
```

```
Out[163]=
```



- Drawing samples from the density--draw from a hat method

- We'll simulate the process of filling a hat with slips of paper, where the number of slips is proportional to the probability the number being in some range (dx1,dx2)

```
In[164]:= m1 = {0,0};
Cov = {{1,.8},{.8,1}};
dx1 = 0.1;
dx2 = dx1;

Nslips=100;
hat = {};
For[x1=-2,x1<=2,x1=x1+dx1,
  For[x2=-2,x2<=2,x2=x2+dx2,
    np = Nslips*multigauss[{x1,x2},m1,Cov];
    For[i=1,i<np,i=i+1,
      hat = Append[hat,{x1,x2}];
    ];
  ];
];
```

hat is a list of pairs of numbers for which the frequency of occurrence of pairs is determined by multigauss.

```
In[171]:= Dimensions[hat]
```

```
Out[171]= {8760, 2}
```

Now let's do a check, where we compile a histogram representing the frequencies of each slip

First, define the "bins" in domain, that we'll use to check for matches:

```
In[172]:= domain = {};
For[x1=-2,x1<=2,x1=x1+dx1,
  For[x2=-2,x2<=2,x2=x2+dx2,
    domain = Append[domain,{x1,x2}];
  ];
];
```

An alternate way of specifying the domain using Outer[], and Range[]:

```
In[174]:= domain2 = Flatten[Outer[List,Range[-2,2,dx1],Range[-2,2,dx1]],1];
```

Now we'll count how many times we find that an element of hat matches a domain element:

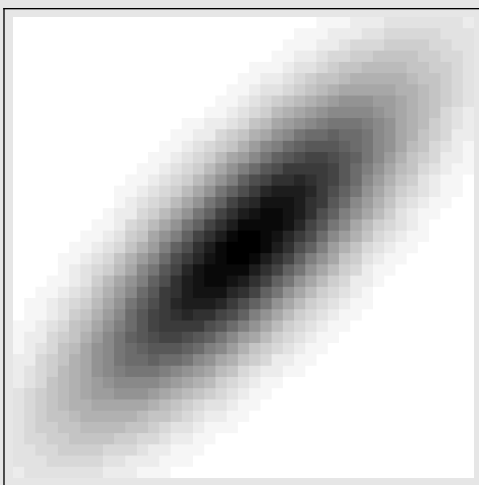
```
In[175]:= Freq = Map[Count[hat,#]&,domain];
```

```
In[176]:= width = Sqrt[Dimensions[Freq]]
```

```
Out[176]= {41}
```

```
In[177]:= ArrayPlot[Partition[Freq,width],ImageSize->Small,DataReversed->True]
```

```
Out[177]=
```



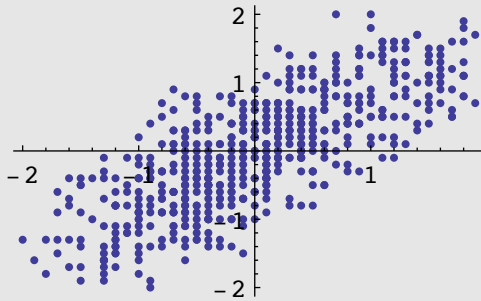
■ Now draw a sample--simulate pulling a slip from hat

```
In[178]:= rv := hat[[RandomInteger[{1,Length[hat]}]]];
```

```
In[179]:= test = Table[hat[[RandomInteger[{1, Length[hat]}]]], {600}];
```

```
In[180]:= g1 = ListPlot[test]
```

```
Out[180]=
```

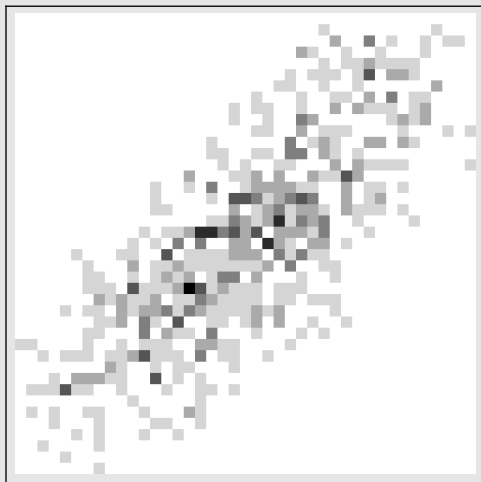


Of course, we can't see the frequency of draws in this plot, so let's count up the number of occurrences per bin, and plot up the results as we did above.

```
In[181]:= Freq2 = Map[Count[test, #] &, domain];  
width = Sqrt[Dimensions[Freq2]];
```

```
In[183]:= ArrayPlot[Partition[Freq2, width], ImageSize -> Small, DataReversed -> True]
```

```
Out[183]=
```



Exercise: Drawing multivariate samples from the density -- use the inverse cumulative distribution

Side comments & where we'll see this again**■ Projection pursuit**

Which projection (marginal) is more "interesting"--the one onto x_1 or onto x_2 ?

Exploratory projection pursuit. (e.g. Intraator, 1993).

■ Inference: Learning parameters of mixture distributions

Return later to the inference problem: Given data, estimate the mixing parameters, means and covariances. EM algorithm.

References

- Applebaum, D. (1996). Probability and Information . Cambridge, UK: Cambridge University Press.
- Cover, T. M., & Joy, A. T. (1991). *Elements of Information Theory*. New York: John Wiley & Sons, Inc.
- Duda, R. O., & Hart, P. E. (1973). Pattern classification and scene analysis . New York.: John Wiley & Sons.
- Golden, R. (1988). A unified framework for connectionist systems. Biological Cybernetics, *59*, 109-120.
- Kersten, D. and P.W. Schrater (2000), *Pattern Inference Theory: A Probabilistic Approach to Vision*, in *Perception and the Physical World*, R. Mausfeld and D. Heyer, Editors. , John Wiley & Sons, Ltd.: Chichester. (pdf)
- Kersten, D., Mamassian P & Yuille A (in press) Object perception as Bayesian inference. Annual Review of Psychology. (pdf, <http://arjournals.annualreviews.org/doi/pdf/10.1146/annurev.psych.55.090902.142005>)
- Kersten, D., & Yuille, A. (2003). Bayesian models of object perception. Current Opinion in Neurobiology, 13(2) <http://gandalf.psych.umn.edu/~kersten/kersten-lab/papers/KerstenYuilleApr2003.pdf>
- Knill, D. C., & Richards, W. (1996). *Perception as Bayesian Inference*. Cambridge: Cambridge University Press.
- Ripley, B. D. (1996). *Pattern Recognition and Neural Networks*. Cambridge, UK: Cambridge University Press.
- Van Trees, H. L. (1968). Detection, Estimation and Modulation Theory . New York: John Wiley and Sons.
- Yuille, A., Coughlan J., Kersten D.(1998) (pdf)