

MATH CLUB

2/28/2013

Note Title

2/28/2013

## THE BEAUTY OF STATISTICAL TESTS

### ① Normal distribution

Random variable  $X$  prob.

Example: Flip a coin  $X = 1$  if heads  $1/2$   
 $X = 0$  if tails  $1/2$

If  $X$  takes on any real value:  
usually, probability  $X =$  given number

$$P\{X = x\} = 0$$

Instead, we have a density function  
= a  $\geq 0$  function on  $\mathbb{R}$ , say,  $f(x)$ .

$$P\{a \leq X \leq b\} = \int_a^b f(x) dx$$

$$\int_{-\infty}^{\infty} f(x) dx = 1. \quad \text{density function} \approx \text{distribution}$$

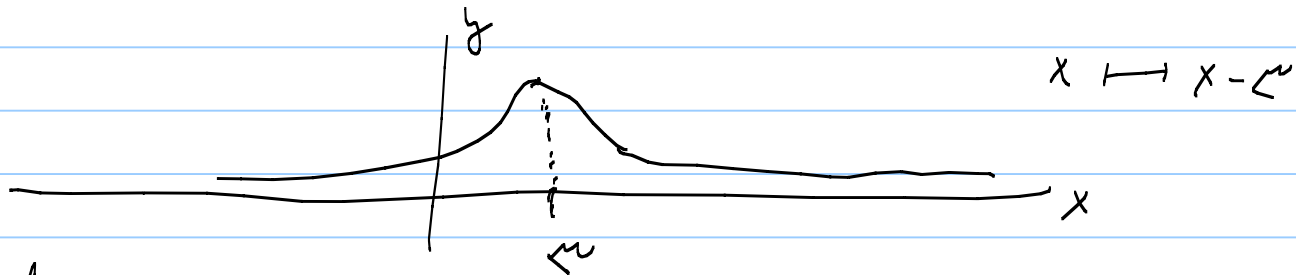
The normal distribution



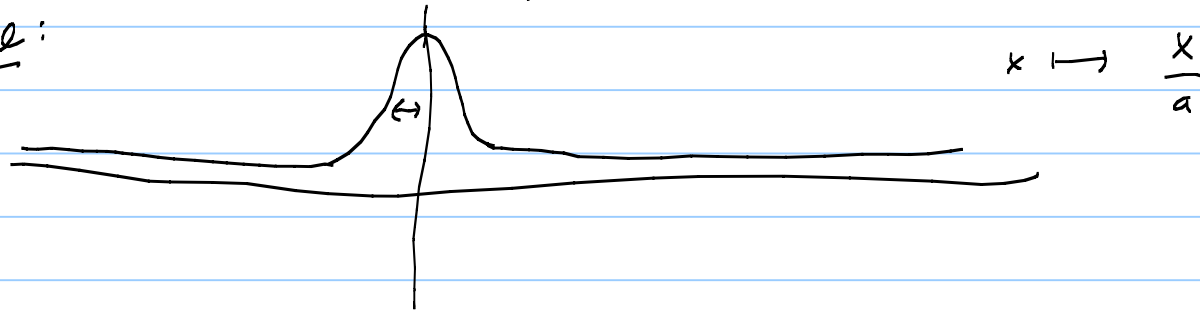
$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

Why is the normal distribution important?

"Obvious transformations": shift



scale:



The more general normal distribution:

$$\frac{1}{a\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2a^2}} \quad \begin{array}{l} \mu = \text{average} \\ = \text{expected value} \\ = E(X) \end{array}$$

$a$  is a measure of the width of the curve

$$a^2 = \text{variance} \quad \text{var}(X) = E((X - E(X))^2).$$

---

The importance of normal distribution:

it does not change (up to changing these two parameters  $\mu, a$ ) when I add independent variables.

Suppose I have two random variables  $X, Y$ .

I can talk about a joint density  $f(x, y)$

$$X = x \quad \& \quad Y = y$$

$$P\{(X, Y) \in S\} = \int_S f(x, y) \, dx \, dy.$$

Independent means  $f(x, y) = g(x)h(y)$

↑ ↑  
individual densities

Key fact: A sum of two independent normal variables is normal.

Application: A sum of a large number of independent random variables with the same distribution, it is approximately normal.

↑  
scale to get finite average and variance.

Informal argument: If such a limit exists,  
the mystery distribution must satisfy the  
"key fact":

$$\underbrace{X_1 + \dots + X_m}_{\sim \text{mystery distribution}} + \underbrace{X_1 + \dots + X_m}_{- \text{mystery distribution}}$$

~ mystery  
distribution

variance prop.

to  $m$

- mystery  
distribution

variance prop.

to  $m$

This uniquely characterizes the normal  
distribution!

Proof sketch of key fact: (think averages = 0)

$$\dots \int e^{-\frac{x^2}{2a^2}} \cdot e^{-\frac{(z-x)^2}{2b^2}} dx$$

examine

dependence on  $z$ .

$$e^{-\frac{(\alpha x + \beta z)^2}{2\gamma^2} - \gamma z^2}$$

bringing to perfect square.



② Multivariate normal distribution



denote

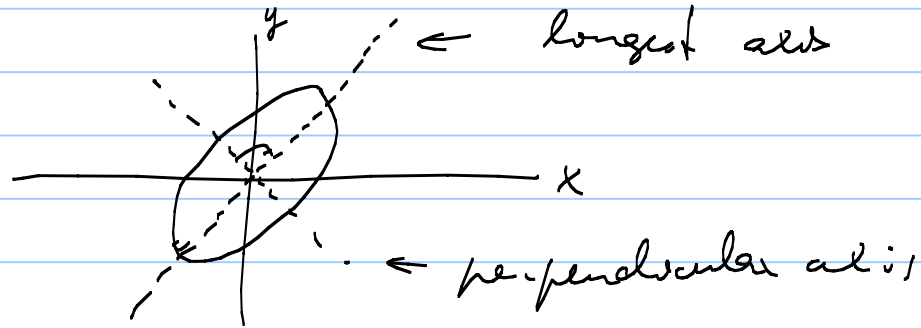
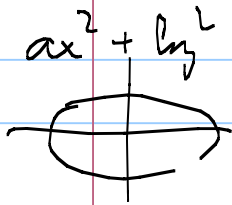
- quadratic expression

negative definite

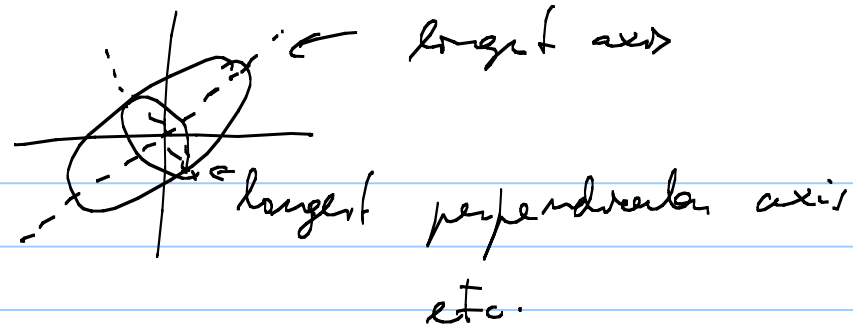
$$\lim_{x \rightarrow \infty} = -\infty$$

two variables: what are the curves  
on which it is constant?

Ellipses:



> 2 variables  
Ellipsoid



→ After turning to perpendicular axes,

The normal distribution

$= ( \quad , \dots , \quad )$  independent  
normal distributions

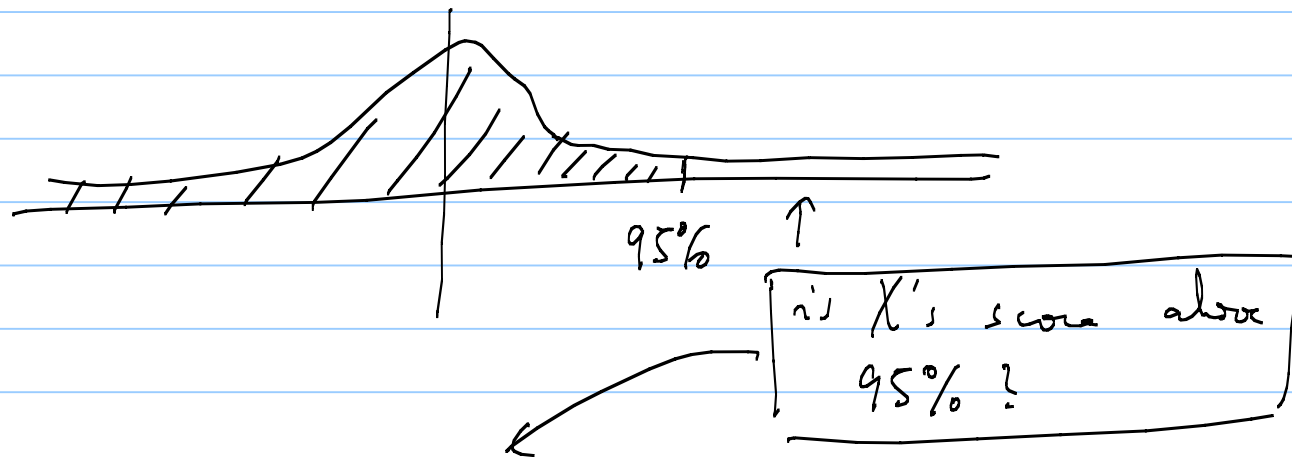
---

③ What is a statistical test?

There is a tutor J. Is this tutor effective?

Take a random student X. Give X a lesson with tutor J. Let X take the test.

Score distribution on the test:



Recipe: If YES  $\rightarrow$  say "Yes, the factor T is effective"

If NO  $\rightarrow$  say nothing.

|              |     | I, X                |       |
|--------------|-----|---------------------|-------|
|              |     | Yes                 | No    |
| T effective? | Yes | Right               | Right |
|              | No  | Wrong<br>$\leq 5\%$ | Right |

$\rightarrow$  null hypothesis

T not effective: 

|    |  |
|----|--|
| 5% |  |
|----|--|

Bayes  
analysis

You can only disprove, but never prove a null hypothesis.

## ④ The Pearson $\chi^2$ test

Draw  $N$  balls with return.  $k$  colors,  $\left. \begin{array}{l} \text{color } i \text{ arises with probability } p_i. \\ \text{hypothesis} \end{array} \right\}$

We get  $a_i$  balls of color  $i$ .

$$\chi^2 = \frac{(a_1 - p_1 N)^2}{p_1 N} + \frac{(a_2 - p_2 N)^2}{p_2 N} + \dots + \frac{(a_k - p_k N)^2}{p_k N}.$$

$\chi^2 \sim \sum_{k=1}^{k-1}$  independent standard normal variables  $\left| \right.$

Example : Cast a die 30 times.



Get 10 times 1

10 times 2

5 times 3

1 times 4

never 5, 6.

Is the die

biased?

$$p_i N = 5$$

$$\chi^2 = \frac{(10-5)^2}{5} + \frac{(10-5)^2}{5} + \frac{(5-5)^2}{5} + \frac{(3-5)^2}{5} + \frac{(0-5)^2}{5} + \frac{(0-5)^2}{5}$$

$$= 20$$

$\chi^2$  with  
 $k-1=5$  degrees  
of freedom

95% 11.07

14.4 > 11.07 Yes!

The doc is 95%  
certainly biased

---

Sketch proof that  $\chi^2$  is distributed as

I said.

①  $\sim$  multivariate normal distribution.  
 $N \gg 0$

Which quadratic expression do I have?

Covariance

$$\text{cov}(X, Y) = E((X - E(X))(Y - E(Y)))$$

Make covariance matrix; look for  
principal axes:

|       |                   |                   |       |                   |
|-------|-------------------|-------------------|-------|-------------------|
|       | $a_1$             | -----             | $a_k$ |                   |
| $a_1$ | $1 - p_1$         | $-\sqrt{p_1 p_2}$ | ----- | $-\sqrt{p_1 p_k}$ |
| ⋮     | $-\sqrt{p_1 p_2}$ | $1 - p_2$         | ----- | $-\sqrt{p_2 p_j}$ |
| ⋮     |                   |                   | ----- |                   |
| $a_k$ |                   |                   | ----- | $1 - p_k$         |

=





then  $\approx \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 0 \end{pmatrix}$



$k-1$  independent standard normal variables!  $\square$