

Mathematical Tools for Neural and Cognitive Science

Fall semester, 2021

Section 4: Summary Statistics & Probability

Statistics is the science of learning from experience, especially experience that arrives a little bit at a time. The earliest information science was statistics, originating in about 1650. This century has seen statistical techniques become the analytic methods of choice in biomedical science, psychology, education, economics, communications theory, sociology, genetic studies, epidemiology, and other areas. Recently, traditional sciences like geology, physics, and astronomy have begun to make increasing use of statistical methods as they focus on areas that demand informational efficiency, such as the study of rare and exotic particles or extremely distant galaxies.

Most people are not natural-born statisticians. Left to our own devices we are not very good at picking out patterns from a sea of noisy data. To put it another way, we are all too good at picking out non-existent patterns that happen to suit our purposes. Statistical theory attacks the problem from both ends. It provides optimal methods for finding a real signal in a noisy background, and also provides strict checks against the overinterpretation of random patterns.

[Efron & Tibshirani, 1998]

Historical context

- 1600's: Early notions of data summary/averaging
- 1700's: Bayesian prob/statistics (Bayes, Laplace)
- 1920's: Frequentist statistics for science (e.g., Fisher)
- 1940's: Statistical signal analysis and communication, estimation/decision theory (e.g., Shannon, Wiener, etc)
- 1950's: Return of Bayesian statistics (e.g., Jeffreys, Wald, Savage, Jaynes...)
- 1970's: Computation, optimization, simulation (e.g., Tukey)
- 2000's: Machine learning (large-scale computing + statistical inference + lots of data)
- Also (since 1950's): statistical neural/cognitive models!

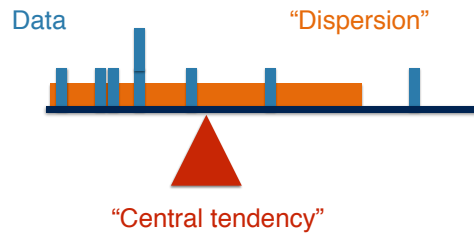
Statistics as summary

0.1, 4.5, -2.3, 0.8, -1.1, 3.2, ...

“The purpose of statistics is to replace a quantity of data by relatively few quantities which shall ... contain as much as possible, ideally the whole, of the relevant information contained in the original data”

- R.A. Fisher, 1934

Descriptive statistics



Descriptive statistics: **Central tendency**

- We often summarize data with averages. Why?
- Average minimizes the squared error (as in regression!):

$$\bar{x} = \arg \min_c \frac{1}{N} \sum_{n=1}^N (x_n - c)^2 = \frac{1}{N} \sum_{n=1}^N x_n = \frac{1}{N} \vec{1}^T \vec{x}$$

- In general: minimize L_p norm: $\arg \min_c \left[\frac{1}{N} \sum_{n=1}^N |x_n - c|^p \right]^{1/p}$
 - $p = 1$: median, m_x
 - $p \rightarrow 0$: mode (location of maximum)
 - $p \rightarrow \infty$: midpoint of range
- Issues: outliers, asymmetry, bimodality

Descriptive statistics: **Dispersion**

- Sample variance (squared standard deviation):

$$s_x^2 = \min_c \frac{1}{N} \sum_{n=1}^N (x_n - c)^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2$$
$$= \frac{1}{N} \sum_{n=1}^N x_n^2 - \bar{x}^2 = \frac{1}{N} \|\vec{x}\|^2 - \bar{x}^2$$

(side note: We'll talk about dividing by N vs. $N-1$ later)

- Mean absolute deviation (MAD) about the median:

$$d_x = \frac{1}{N} \sum_{n=1}^N |x_n - m_x|$$

- Quantiles (eg: "90% of data lie in range [1.5 8.2]")

A note on notation

- We have, and will continue to use the notation for a "sample mean" (\bar{x}) and a "sample standard deviation" (s) or variance (s^2).
- Statistics makes a distinction between these sample values and the corresponding "population" values of mean (μ) and variance (σ^2).
- We'll return to this distinction later on in the for now only consider sample data and stick corresponding notation.

This really doesn't belong here! ... we haven't introduced probability yet, and so talking about μ and σ^2 doesn't make sense. Should wait until after we introduce probability, and then want to get back to data/stats....

Descriptive statistics: Multi-D

Data points: $\{\vec{d}_n\}$ $n \in [1 \dots N]$, in 2-D: $\vec{d}_n = \begin{bmatrix} x_n \\ y_n \end{bmatrix}$

As in 1D: define central tendency as vector that minimizes the sum of squared distances to all data points:

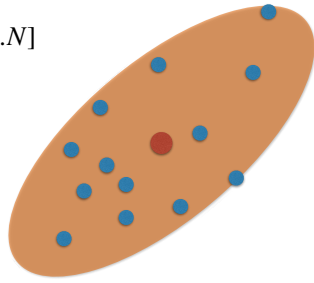
$$\vec{d} \equiv \arg \min_{\vec{c}} \sum_n \|\vec{d}_n - \vec{c}\|^2$$
$$= \arg \min_{c_x, c_y} \sum_n (x_n - c_x)^2 + (y_n - c_y)^2 \quad (\text{in 2-D})$$
$$= \frac{1}{N} \sum_{n=1}^N \vec{d}_n = \begin{bmatrix} \bar{x} \\ \bar{y} \end{bmatrix}$$

Descriptive statistics: Multi-D

Data points: $\{\vec{d}_n\} \quad n \in [1 \dots N]$

Sample mean (average):

$$\bar{\vec{d}} = \frac{1}{N} \sum_{n=1}^N \vec{d}_n = \begin{bmatrix} \bar{x} \\ \bar{y} \end{bmatrix}$$



Sample covariance:

$$\begin{aligned} C_d &= \frac{1}{N} \sum_{n=1}^N (\vec{d}_n - \bar{\vec{d}})(\vec{d}_n - \bar{\vec{d}})^T = \frac{1}{N} \sum_{n=1}^N \vec{d}_n \vec{d}_n^T - \bar{\vec{d}} \bar{\vec{d}}^T \\ &= \frac{1}{N} \begin{bmatrix} ||\vec{x}||^2 & \vec{x}^T \vec{y} \\ \vec{y}^T \vec{x} & ||\vec{y}||^2 \end{bmatrix} - \begin{bmatrix} \bar{x}^2 & \bar{x} \bar{y} \\ \bar{y} \bar{x} & \bar{y}^2 \end{bmatrix} = \begin{bmatrix} s_x^2 & s_{xy} \\ s_{xy} & s_y^2 \end{bmatrix} \end{aligned}$$

Affine transformations

If $\vec{b}_n = M(\vec{d}_n - \vec{a})$ (translate, then rotate-stretch-rotate)

then $\bar{\vec{b}} = M(\bar{\vec{d}} - \vec{a})$

$$C_b = M C_d M^T$$

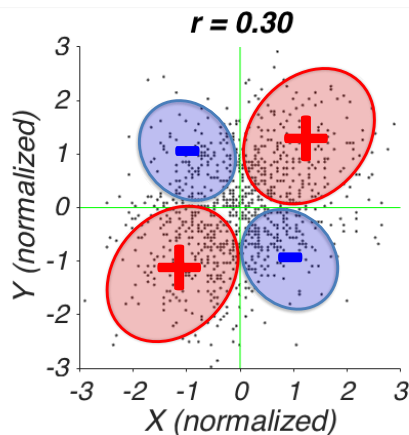
Special case: “re-center” and “normalize” the data:

$$\vec{a} = \begin{bmatrix} \bar{x} \\ \bar{y} \end{bmatrix} \quad M = \begin{bmatrix} \frac{1}{s_x} & 0 \\ 0 & \frac{1}{s_y} \end{bmatrix}$$

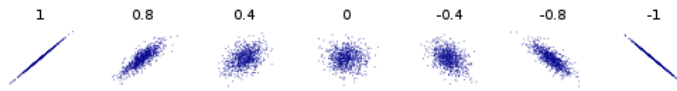
$$\text{then } \bar{\vec{b}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad C_b = \begin{bmatrix} 1 & \frac{s_{xy}}{s_x s_y} \\ \frac{s_{xy}}{s_x s_y} & 1 \end{bmatrix} \quad [on \ board]$$

“,”
(Pearson
correlation
coefficient)

Correlation



Correlation r captures dependency



... but not slope!



Regression (revisited)

$$\vec{y} = \beta \vec{x} + \vec{e}$$

Optimal regression line slope:

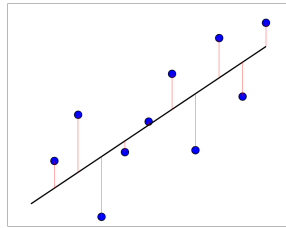
$$\beta = \frac{\vec{x}^T \vec{y}}{\vec{x}^T \vec{x}} = \frac{s_{xy}}{s_x^2}$$

Error variance:

$$s_e^2 = s_y^2 - 2\beta s_{xy} + \beta^2 s_x^2$$

$$= s_y^2 - \frac{s_{xy}^2}{s_x^2}$$

Partition of variance:
error variance = data variance - explained variance



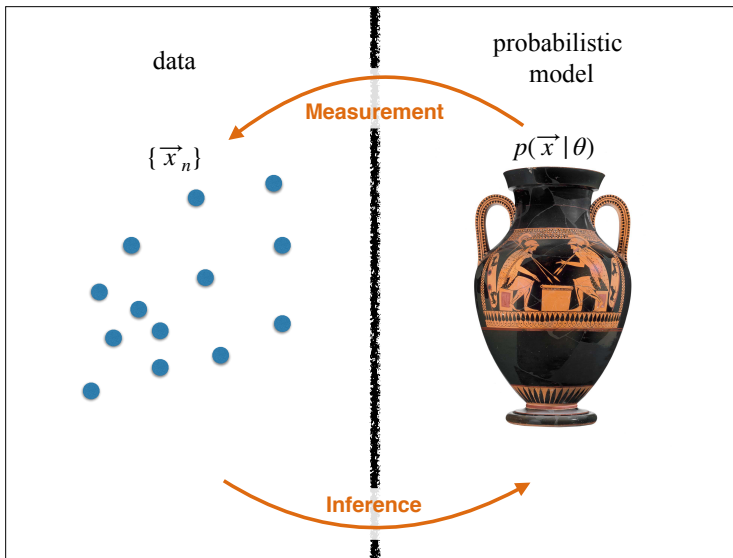
Expressed as a proportion of σ_y^2 :

$$\frac{s_e^2}{s_y^2} = 1 - \frac{s_{xy}^2}{s_x^2 s_y^2} = 1 - r^2$$

“r-squared”
(proportion
of variance
explained)

Probability: an abstract mathematical framework for describing random quantities, or stochastic models of the world

Statistics: use of probability to summarize, analyze, interpret data. **Fundamental to all experimental science.**



Probabilistic Middleville

In Middleville, every family has two children, brought by the stork.

The stork delivers boys and girls randomly, with family probabilities $\{BB, BG, GB, GG\} = \{0.2, 0.3, 0.2, 0.3\}$

probabilistic model

You pick a family at random and discover that one of the children is a girl.

new data

What are the chances that the other child is a girl?

inference

Statistical Middleville

In Middleville, every family has two children, brought by the stork.

In a survey of 100 of the Middleville families, 32 have two girls, 23 have two boys, and the remainder one of each.

prev data

You pick a family at random and discover that one of the children is a girl.

new data

What are the chances that the other child is a girl?

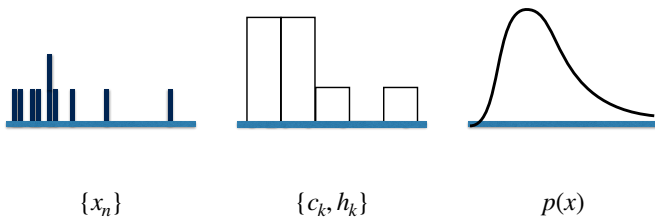
inference

Univariate Probability (outline)

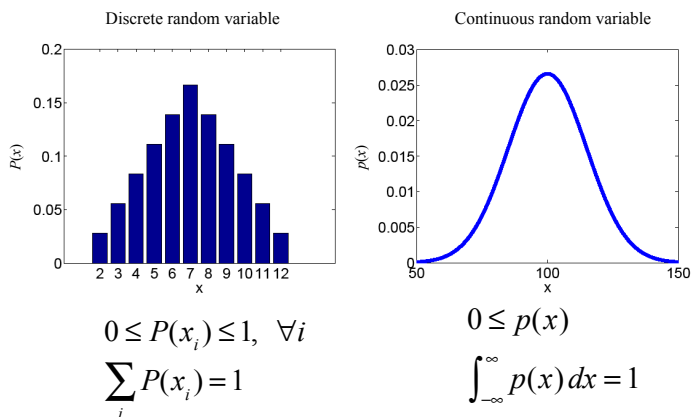
- distributions: discrete and continuous
- expected value, moments
- transformations: affine, monotonic nonlinear
- cumulative distributions. Quantiles, drawing samples

Frequentist view of probability: limit of infinite data

data \rightarrow histogram \rightarrow probability distribution

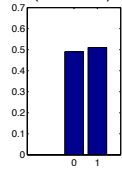


Probability distributions

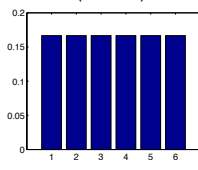


Example distributions

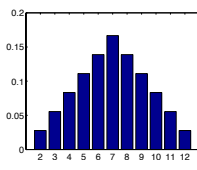
a not-quite-fair coin
(Bernoulli)



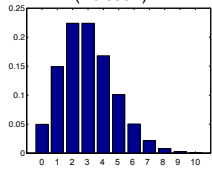
roll of a fair die
(uniform)



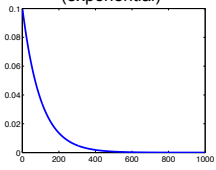
sum of rolls of
two fair dice



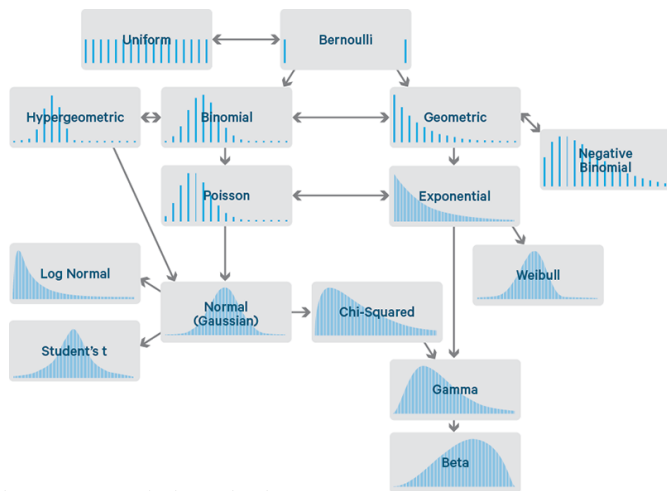
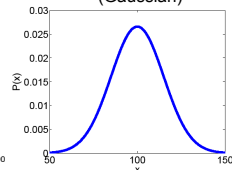
clicks of a Geiger counter,
in a fixed time interval
(Poisson)



... and, time between clicks
(exponential)



horizontal velocity of gas
molecules exiting a fan
(Gaussian)



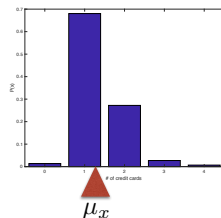
[Figure: Sean Owen, Cloudera Engineering]

Expected value (discrete random variable)

$$\mathbb{E}(x) = \sum_{k=1}^K x_k p(x_k)$$

(the mean, μ_x)

a weighted sum over the discrete values



More generally: $\mathbb{E}(f(x)) = \sum_{k=1}^K f(x_k) p(x_k)$ (sum over discrete values)

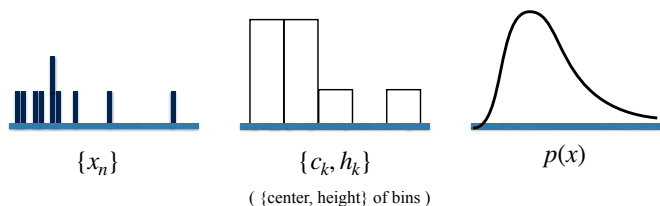
Sample average, an estimate of the expected value:

$$\mathbb{E}(f(x)) \approx \bar{f}(x) = \frac{1}{N} \sum_{n=1}^N f(x_n) \quad (\text{sum is over data samples!})$$

Sample average converges to expected value as one gathers more data...

Expected value (continuous random variable)

data → histogram → probability distribution



$$\bar{x} = \frac{1}{N} \sum_n x_n \quad \bar{x} \approx \frac{1}{K} \sum_k c_k h_k = \vec{c}^T \vec{h} \quad \mu_x = \int x p(x) dx$$

Expected value (continuous)

$$\mathbb{E}(x) = \int x p(x) dx \quad [\text{“mean”, } \mu]$$

$$\mathbb{E}(x^2) = \int x^2 p(x) dx \quad [\text{“second moment”, } m_2]$$

$$\begin{aligned} \mathbb{E}((x - \mu)^2) &= \int (x - \mu)^2 p(x) dx & [\text{“variance”, } \sigma^2] \\ &= \int x^2 p(x) dx - \mu^2 & [m_2 \text{ minus } \mu^2] \end{aligned}$$

$$\mathbb{E}(f(x)) = \int f(x) p(x) dx \quad [\text{“expected value of } f\text{”}]$$

Note: expectation is an inner product, and thus *linear*, so:

$$\mathbb{E}(af(x) + bg(x)) = a\mathbb{E}(f(x)) + b\mathbb{E}(g(x))$$

Transformations of scalar random variables

$Y = aX + b$ “affine” (linear plus constant)

Analogous to sample mean/covariance:

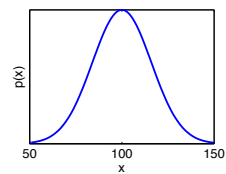
$$\mu_Y = \mathbb{E}(Y) = a\mathbb{E}(X) + b = a\mu_X + b$$

$$\sigma_Y^2 = \mathbb{E}((Y - \mu_Y)^2) = \mathbb{E}((aX - a\mu_X)^2) = a^2 \sigma_X^2$$

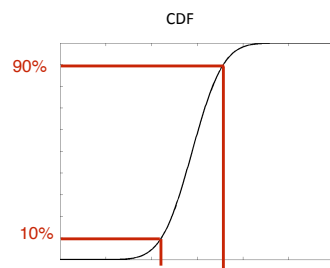
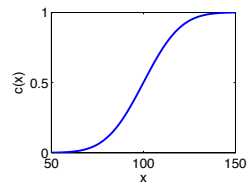
Full distribution: $p_Y(y) = \frac{1}{a} p_X\left(\frac{y - b}{a}\right)$

$Y = g(X)$ “monotonic” (derivative > 0)

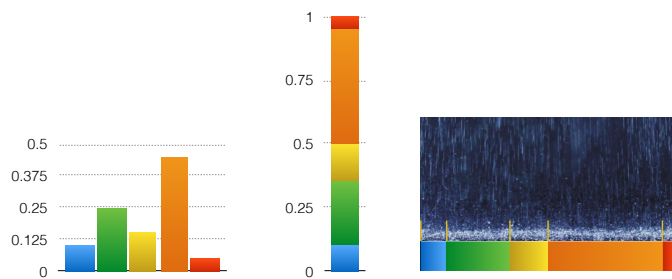
$$p_Y(y) = \frac{p_X(g^{-1}(y))}{g'(g^{-1}(y))}$$



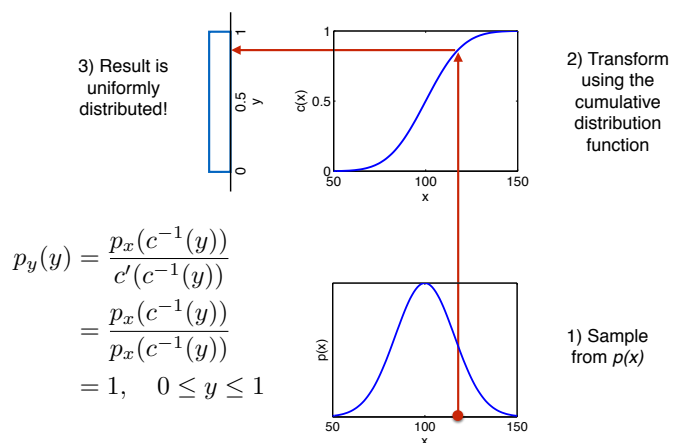
$$c(x) = \int_{-\infty}^x p(z) \, dz$$



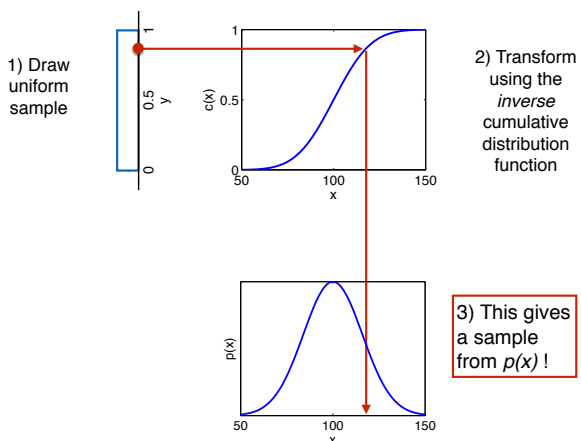
Drawing samples - discrete



Drawing samples - continuous



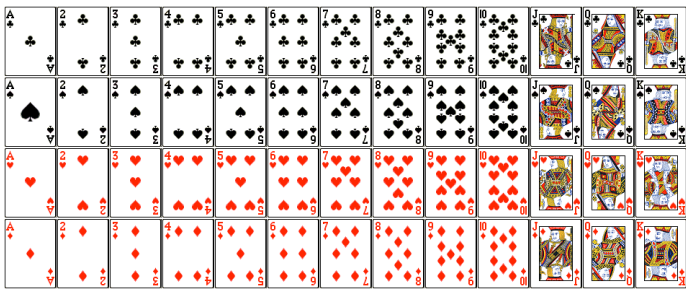
Drawing samples - continuous



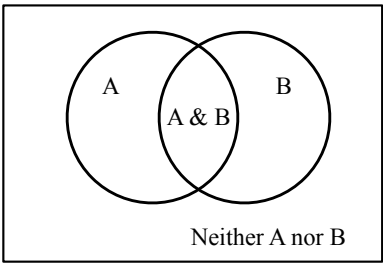
Multi-variate probability (outline)

- Joint distributions
- Marginals (integrating)
- Conditionals (slicing)
- Bayes’ rule (inverse probability)
- Statistical independence (separability)
- Mean/Covariance
- Linear transformations

Joint and conditional probability - discrete

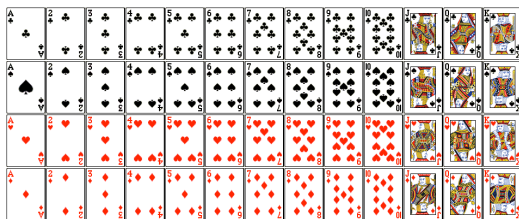


Conditional probability



$p(A|B)$ = probability of A given that B is asserted to be true = $\frac{p(A \& B)}{p(B)}$

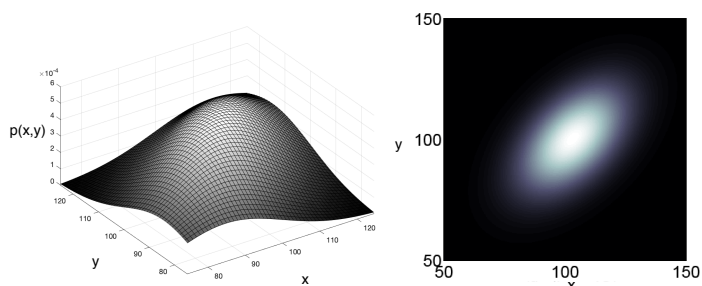
Joint and conditional probability - discrete



$P(\text{Ace})$
 $P(\text{Heart})$
 $P(\text{Ace \& Heart})$
 $P(\text{Ace} \mid \text{Heart})$
 $P(\text{not Jack of Diamonds})$
 $P(\text{Ace} \mid \text{not Jack of Diamonds})$

“Independence”

Joint distribution (continuous)

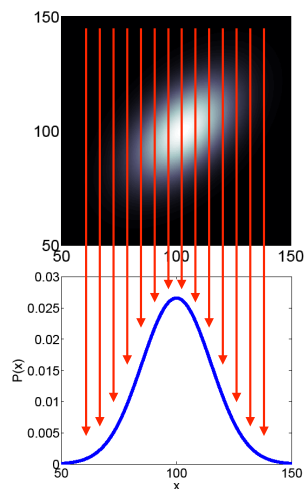


$p(x, y)$

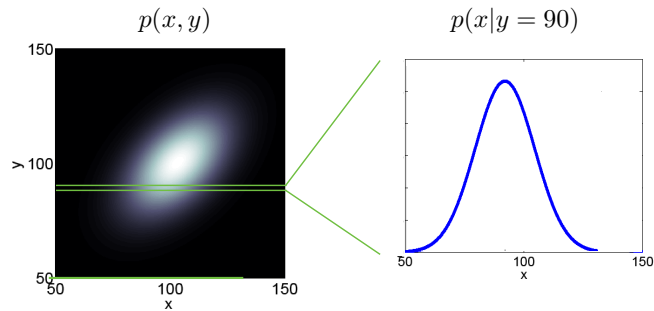
Marginal distribution

$p(x, y)$

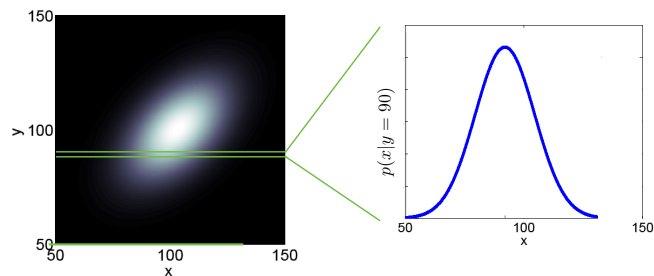
$$p(x) = \int p(x, y) dy$$



Conditional distribution



Conditional distribution



$$p(x|y = 90) = p(x, y = 90) / \int p(x, y = 90) dx$$

$$= \underbrace{p(x, y = 90)}_{\text{slice joint distribution}} / \underbrace{p(y = 90)}_{\text{normalize (by marginal)}}$$

More generally:

$$p(x|y) = p(x, y) / p(y)$$

slice joint distribution

normalize (by marginal)

Bayes' Rule



LII. *An Essay towards solving a Problem in the Doctrine of Chances.* By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a Letter to John Canton, A. M. F. R. S.

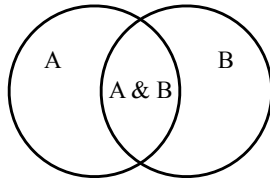
Dear Sir,

Read Dec. 23, 1763. **I** Now send you an essay which I have found among the papers of our deceased friend Mr. Bayes, and which, in my opinion, has great merit, and well deserves to be preserved.

$$p(x|y) = p(y|x) p(x) / p(y)$$

(a direct consequence of the definition of conditional probability)

Bayes' Rule



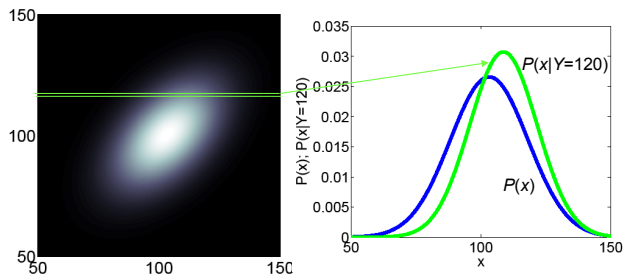
$p(A|B)$ = probability of A given that B is asserted to be true = $\frac{p(A \& B)}{p(B)}$

$$p(A \& B) = p(B)p(A|B)$$

$$= p(A)p(B|A)$$

$$\Rightarrow p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Conditional vs. marginal



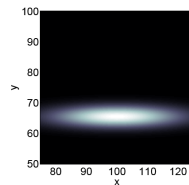
In general, the marginals for different Y values differ.

When are they the same? In particular, when are all conditionals equal to the marginal?

Statistical independence

Random variables X and Y are statistically independent if (and only if):

$$p(x, y) = p(x)p(y) \quad \forall x, y$$



(note: for discrete distributions, this is an outer product!)

Independence implies that *all* conditionals are equal to the corresponding marginal:

$$p(x|y) = p(x, y) / p(y) = p(x) \quad \forall x, y$$

Mean, covariance, affine transformations

For R.V. \vec{x} , $\vec{\mu}_x = \mathbb{E}(\vec{x})$, $C_x = \mathbb{E}((\vec{x} - \vec{\mu}_x)(\vec{x} - \vec{\mu}_x)^T)$

For R.V. $\vec{y} = M(\vec{x} - \vec{a})$,

analogous to results for sample mean/covariance:

$$\vec{\mu}_y = \mathbb{E}(M(\vec{x} - \vec{a}))$$

$$= M(\mathbb{E}(\vec{x}) - \vec{a})$$

$$= M(\vec{\mu}_x - \vec{a})$$

$$C_y = \mathbb{E}((M(\vec{x} - \vec{\mu}_x))(M(\vec{x} - \vec{\mu}_x))^T)$$

$$= M\mathbb{E}((\vec{x} - \vec{\mu}_x)(\vec{x} - \vec{\mu}_x)^T)M^T$$

$$= MC_x M^T$$

Special case: Sum of two RVs

Let $Z = X + Y$, or $Z = \vec{1}^T \begin{bmatrix} X \\ Y \end{bmatrix}$

$$\mu_Z = \mu_X + \mu_Y$$

$$\sigma_Z^2 = \sigma_X^2 + 2\sigma_{XY} + \sigma_Y^2$$

Special case: if X and Y are *independent*, then:

$$\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y) \text{ and thus } \sigma_{XY} = 0$$

$$\sigma_Z^2 = \sigma_X^2 + \sigma_Y^2$$

$p_Z(z)$ is the *convolution* of $p_X(x)$ and $p_Y(y)$

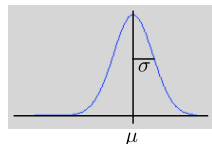
[on board]

Gaussian (a.k.a. "Normal") densities

One-dimensional:

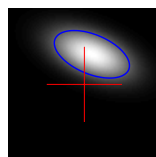
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Alt. notation: $x \sim N(\mu, \sigma^2)$



Multi-dimensional:

$$p(\vec{x}) = \frac{1}{\sqrt{(2\pi)^N |C|}} e^{-(\vec{x} - \vec{\mu})^T C^{-1} (\vec{x} - \vec{\mu}) / 2}$$



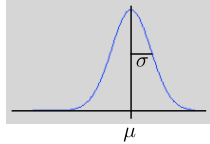
mean: [0.2, 0.8]

cov: [1.0 -0.3;
-0.3 0.4]

Gaussian properties

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

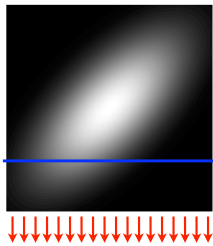
$$p(\vec{x}) = \frac{1}{\sqrt{(2\pi)^N |C|}} e^{-\frac{(\vec{x}-\vec{\mu})^T C^{-1} (\vec{x}-\vec{\mu})}{2}}$$



- joint density of indep Gaussian RVs is elliptical [easy]
- conditionals of a Gaussian are Gaussian [easy]
- marginals of a Gaussian are Gaussian [easy]
- product of two Gaussian dists is Gaussian [easy]
- sum of independent Gaussian RVs is Gaussian [moderate]
- the most random (max entropy) density of given variance [moderate]
- central limit theorem: sum of many indep. RVs is Gaussian [hard]

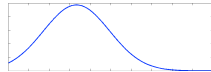
let $P = C^{-1}$ (the “precision” matrix)

$$\begin{aligned} p(x_1|x_2=a) &\propto e^{-\frac{1}{2}[P_{11}(x_1-\mu_1)^2+2P_{12}(x_1-\mu_1)(a-\mu_2)+\dots]} \\ &= e^{-\frac{1}{2}[P_{11}x_1^2+2(P_{12}(a-\mu_2)-P_{11}\mu_1)x_1+\dots]} \\ &= e^{-\frac{1}{2}\left(x_1-\mu_1+\frac{P_{12}}{P_{11}}(a-\mu_2)\right)P_{11}\left(x_1-\mu_1+\frac{P_{12}}{P_{11}}(a-\mu_2)\right)+\dots} \end{aligned}$$

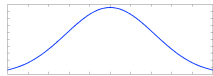


Gaussian, with: $\mu = \mu_1 - \frac{P_{12}}{P_{11}}(a - \mu_2)$
 $\sigma^2 = \frac{1}{P_{11}}$

Conditional:



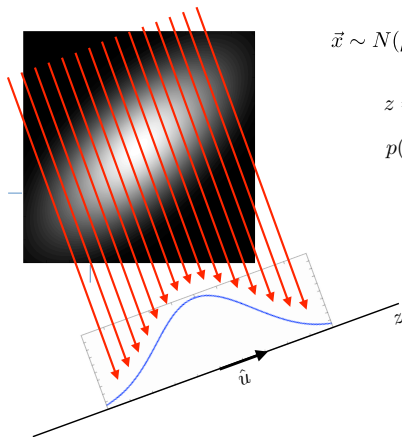
Marginal:



$$p(x_1) = \int p(\vec{x}) dx_2 \quad [on\ board]$$

Gaussian, with: $\mu = \mu_1$
 $\sigma^2 = C_{11}$

Generalized marginals of a Gaussian



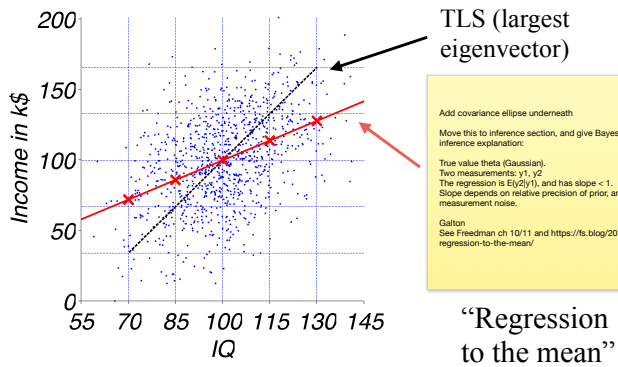
$$\vec{x} \sim N(\vec{\mu}_x, C_x)$$

$$z = \hat{u}^T \vec{x}$$

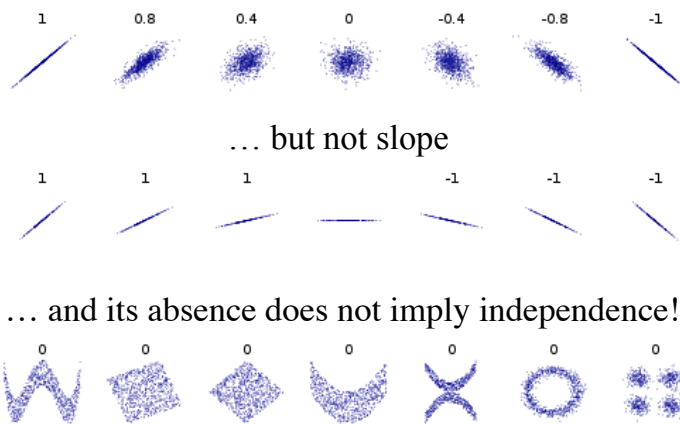
$p(z)$ is Gaussian, with:

$$\begin{aligned} \mu_z &= \hat{u}^T \vec{\mu}_x \\ \sigma_z^2 &= \hat{u}^T C_x \hat{u} \end{aligned}$$

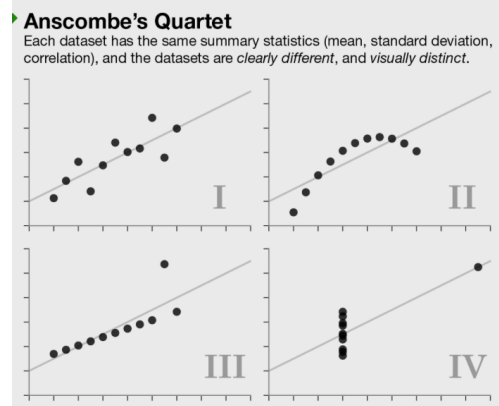
Correlation and regression

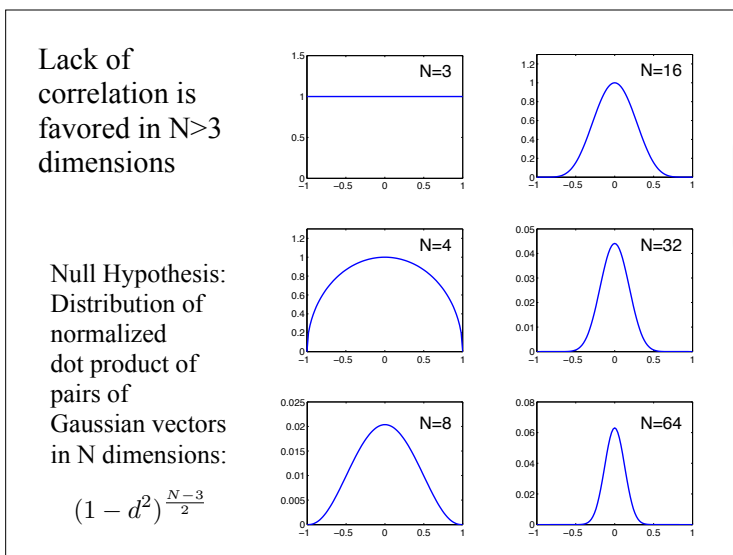
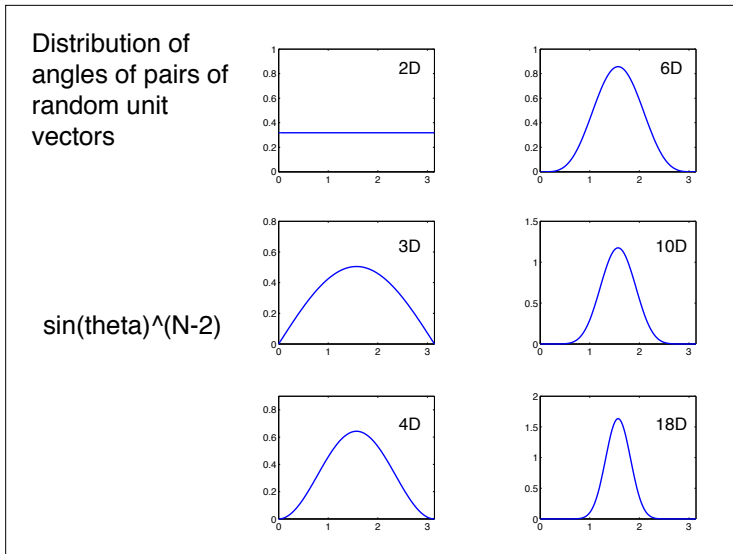
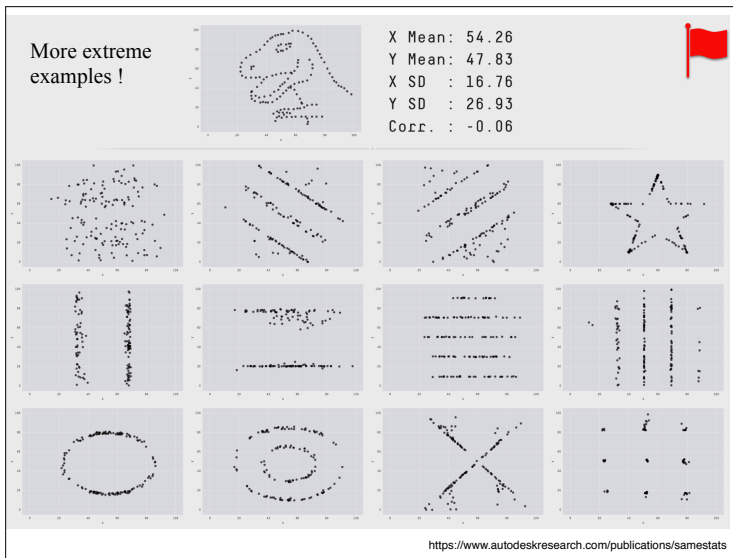


Correlation implies dependency

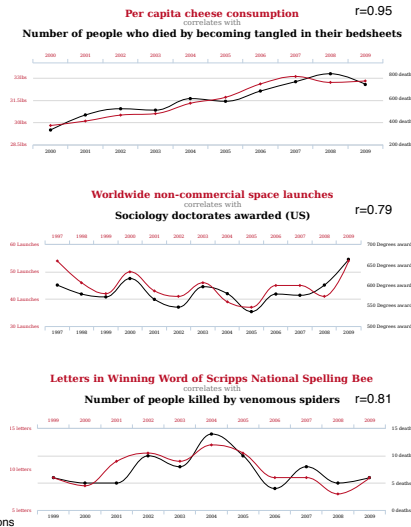


Correlation between variables does not uniquely indicate the shape of their joint distribution





Nevertheless,
one can find
correlation if
one looks for it!

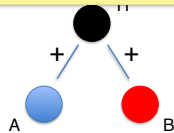


Covariation/correlation does not imply causation

- Correlation does not provide a direction for causality. For that, you need additional (temporal) information.
- More generally, correlations are often a result of hidden (unmeasured, uncontrolled) variables

Example: conditional independence:

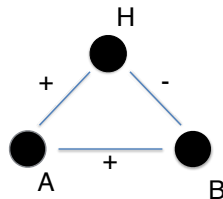
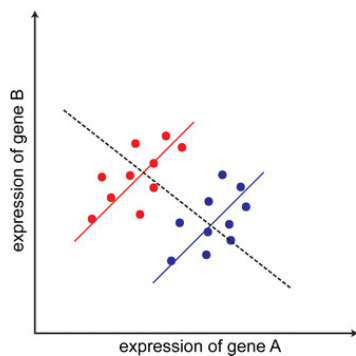
$$p(A, B | H) = p(A | H)p(B | H)$$



Move these next few slides to inference section, after defining Gaussian MAP estimation.
Then, this example will make more sense, and can be

[On board: in Gaussian case, connections are explicit in the precision matrix]

Another example: "Simpson's paradox"

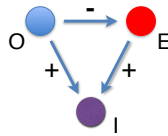


Milton Friedman's Thermostat



O = outside temperature (assumed cold)
I = inside temperature (ideally, constant)
E = energy used for heating

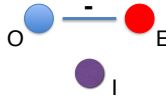
True interactions:



Statistical observations:

- O and I uncorrelated
- I and E uncorrelated
- O and E anti-correlated

Statistical interactions, $P=C^{-1}$:



Some nonsensical conclusions:

- O and E have no effect on I, so shut off heater to save money!
- I is irrelevant, and can be ignored. Increases in E cause decreases in O.

Statistical summary cannot replace scientific reasoning/experiments!

Summary: Correlation misinterpretations



- Correlation implies dependency, but lack of correlation does *not* imply independence
- Correlation does *not* imply data lie near a line/plane/hyperplane (subspace), with simple noise perturbations
- Correlation does *not* imply causation (temporally, or by direct influence/connection)
- Correlation is a descriptive statistic, and does not eliminate the need for scientific reasoning/experiment!