# Mathematical Tools for Neural and Cognitive Science

Fall semester, 2025

# 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]

#### Some 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 (statistical inference with large-scale computing + lots of data)
- Also (since 1950's): statistical neural/cognitive models!

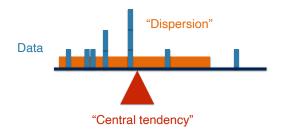
#### Statistics as summary description

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

#### Standard descriptive statistics



#### Descriptive statistics: 1D

The most common measures of central tendency & dispersion:

• Sample mean minimizes the squared error

$$m_d = \arg\min_{c} \frac{1}{N} \sum_{n=1}^{N} (d_n - c)^2$$
$$= \frac{1}{N} \sum_{n=1}^{N} d_n$$

 $s_d^2$ 

• Sample variance is the squared error

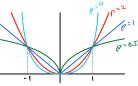
$$\begin{split} s_d^2 &= \min_c \frac{1}{N} \sum_{n=1}^N \left(d_n - c\right)^2 = \frac{1}{N} \sum_n \left(d_n - m_d\right)^2 \\ &= \frac{1}{N} \sum_n d_n^2 - m_d^2 \qquad \qquad \text{(second moment minus squared mean)} \end{split}$$

#### Descriptive statistics: generalizations

More generally, can measure dispersion with

an " $L_p$  norm":

$$\left[\sum_{n=1}^{N} \left| d_n - c \right|^p \right]^{1/p}$$



Different *p* values give different measures of **central tendency**:

- p=2 : mean (standard choice)
- p = 1 : median
- $p \to 0$  : mode (location of maximum)
- $p \to \infty$  : midpoint of range

#### Descriptive statistics: 2-D

- Data points:  $\vec{d_n} = \begin{bmatrix} x_n \\ y_n \end{bmatrix}, n \in [1 \dots N]$
- Sample mean: the vector that minimizes average squared distance to data points:

$$ec{m}_d = \arg\min_{ec{c}} rac{1}{N} \sum_{n=1}^N \| ec{d}_n - ec{c} \, \|^2, \quad ec{c} = egin{bmatrix} c_x \ c_y \end{bmatrix}$$

$$= \arg\min_{c_x, c_y} \frac{1}{N} \sum_{n} \left[ (x_n - c_x)^2 + (y_n - c_y)^2 \right]$$

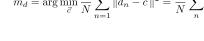
$$= \frac{1}{N} \sum_{n} \begin{bmatrix} x_n \\ y_n \end{bmatrix} = \frac{1}{N} \sum_{n} \vec{d_n}$$

(analogous to scalar case!)

# Descriptive statistics: 2-D

- $\bullet \ \ \text{ Data points: } \ \vec{d_n} = \begin{bmatrix} x_n \\ y_n \end{bmatrix}, \quad n \in [1 \dots N]$
- Sample mean:

$$\vec{m}_d = \arg\min_{\vec{c}} \frac{1}{N} \sum_{n=1}^{N} ||\vec{d}_n - \vec{c}||^2 = \frac{1}{N} \sum_{n} \vec{d}_n$$



O Sample (total) variance:

$$s_d^2 = \min_{\vec{c}} \frac{1}{N} \sum_{n=1}^N ||\vec{d}_n - \vec{c}||^2 = \frac{1}{N} \sum_n ||\vec{d}_n - \vec{m}_d||^2$$

$$= \frac{1}{N} \sum_n ||\vec{d}_n||^2 - ||\vec{m}_d||^2$$
 (analogous)

$$= \frac{1}{N} \sum_{n} \|\vec{d}_n\|^2 - \|\vec{m}_d\|^2$$

#### Descriptive statistics: 2-D

• Sample mean, in direction  $\hat{u}$ 

$$m_u = \arg\min_c \frac{1}{N} \sum_n \left( \hat{u}^T \vec{d}_n - c \right)^2$$
$$= \frac{1}{N} \sum_n \hat{u}^T \vec{d}_n = \hat{u}^T \frac{1}{N} \sum_n \vec{d}_n = \hat{u}^T \vec{m}_d$$

• Sample variance, in direction  $\,\hat{u}\,$ 

Sample variance, in direction 
$$u$$

$$s_u^2 = \min_c \frac{1}{N} \sum_{n=1}^N \left( \hat{u}^T \vec{d}_n - c \right)^2 = \frac{1}{N} \sum_n \left( \hat{u}^T \vec{d}_n - \hat{u}^T \vec{m}_d \right)^2$$

$$= \hat{u}^T \left[ \frac{1}{N} \sum_n \left( \vec{d}_n - \vec{m}_d \right) \left( \vec{d}_n - \vec{m}_d \right)^T \right] \hat{u}$$

$$= \hat{u}^T \left[ \frac{1}{N} \sum_n \left( \vec{d}_n \vec{d}_n^T - \vec{m}_d \vec{m}_d^T \right) \right] \hat{u}$$
sample covariance,  $C_d$ 

#### Descriptive statistics: 2-D

• Sample mean, in direction  $\,\hat{u}\,$ 

$$\begin{split} m_u &= \arg\min_c \frac{1}{N} \sum_n \left( \hat{u}^T \vec{d_n} - c \right)^2 \\ &= \frac{1}{N} \sum_n \hat{u}^T \vec{d_n} = \hat{u}^T \frac{1}{N} \sum_n \vec{d_n} = \hat{u}^T \vec{m}_d \end{split}$$

• Sample variance, in direction  $\hat{u}$ 

$$s_{u}^{2} = \hat{u}^{T} \begin{bmatrix} \frac{1}{N} \sum_{n} \left( \vec{d}_{n} \vec{d}_{n}^{T} - \vec{m}_{d} \vec{m}_{d}^{T} \right) \end{bmatrix} \hat{u}$$

$$\frac{1}{N} \sum_{n} \begin{bmatrix} x_{n}^{2} & x_{n} y_{n} \\ y_{n} x_{n} & x_{n}^{2} \end{bmatrix} - \begin{bmatrix} m_{x}^{2} & m_{x} m_{y} \\ m_{y} m_{x} & m_{y}^{2} \end{bmatrix} = \begin{bmatrix} s_{x}^{2} & s_{xy} \\ s_{xy} & s_{y}^{2} \end{bmatrix}$$

# Descriptive statistics: multi-D

• Data points: matrix D (N data vectors in *columns*)

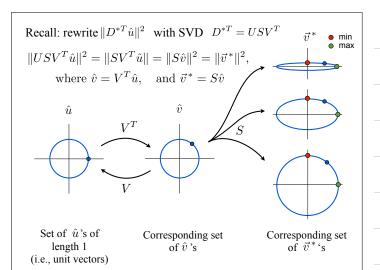
Sample mean:

$$\vec{m}_d = \frac{1}{N} \sum_n \vec{d}_n = \frac{1}{N} D \vec{1}$$

• Sample variance, in direction  $\hat{u}$ 

$$\begin{split} s_u^2 &= \frac{1}{N} \sum_n \left( (\vec{d_n} - \vec{m}_d)^T \hat{u} \right)^2 = \frac{1}{N} || \underbrace{\left( D - \vec{m}_d \vec{\mathbf{1}}^T \right)^T}_{D^*} \hat{u} ||^2 \\ &= \frac{1}{N} || D^{*T} \hat{u} ||^2 \end{split}$$

As we vary direction  $\hat{u}$ , what does the sample variance do?



#### Descriptive statistics: multi-D

- Data points: matrix D (N data vectors in *columns*)
- **Sample mean, in direction**  $\hat{u}$ :

$$m_u = \frac{1}{N}(\hat{u}^T D)\vec{1} = \hat{u}^T \boxed{\begin{bmatrix} D \ \vec{1} \\ N \end{bmatrix}}$$

sample mean,  $\vec{m}_d$ 

O Sample variance, in direction û:

$$\begin{split} s_u^2 &= \frac{1}{N} \sum_n \left( (\vec{d_n} - \vec{m_d})^T \hat{u} \right)^2 = \frac{1}{N} \| (D - \vec{m_d} \vec{\mathbf{1}}^T)^T \hat{u} \|^2 \\ &= \frac{1}{N} \| D^{*T} \hat{u} \|^2 = \hat{u}^T \boxed{\frac{D^* D^{*T}}{N}} \hat{u} \\ &\text{sample covariance, } C_d \end{split}$$

# Principal Component Analysis (PCA)

The shape of a data cloud can be summarized with an ellipse (ellipsoid), centered around the mean, using a simple procedure:

- (1) Subtract mean from all data points (re-centers data around origin)
- (2) Collect centered data vectors in columns of a matrix,  $D^*$
- (3) Compute the SVD:  $D^{*T} = USV^T$

or use covariance matrix  $C_d = D^*D^{*T} = V S^T S^T S^T$ 

 $\Lambda$ , square and diagonal, elements  $\lambda_k$ 

- Columns of V are the *principal components* (axes) of the ellipsoid, singular values  $s_k$  (or  $\sqrt{\lambda_k}$ ) are the corresponding *principal radii*.
- Ellipse volume is proportional to product of  $s_k$ 's.
- Total variance is equal to sum of  $\lambda_k$ 's.



#### Eigenvectors/eigenvalues

- An *eigenvector* of a matrix is a vector that is rescaled by the matrix (i.e., the direction is unchanged)
- The corresponding scale factor is called the *eigenvalue*
- For covariance matrix  $C_d = D^*D^{*T} = V\Lambda V^T$  the columns of V (denoted  $\hat{v}_k$ ) are eigenvectors, with corresponding eigenvalues  $\lambda_k$ :

$$C_d \hat{v}_k = V \Lambda V^T \hat{v}_k$$

$$= V \Lambda \hat{e}_k$$

$$= \lambda_k V \hat{e}_k$$

$$= \lambda_k \hat{v}_k$$

 $\bullet$  For LSI system L, the eigenvectors are complex exponentials:

$$L\vec{v}_k = FRF^T\vec{v}_k = r_k\vec{v}_k$$

F is the Fourier transform,  $\vec{v}_k$  the kth Fourier basis function,  $r_k$  the kth entry of diagonal matrix R containing F.T. of impulse response

#### Affine transformations

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

then 
$$\vec{m}_b = M \left( \vec{m}_d - \vec{a} \right)$$
 (mean and covariance transform according to simple rules)

$$C_b = MC_d M^T$$

Standard case: "re-center" and "normalize" the components:

Let 
$$\vec{a} = \vec{m}_d$$
  $M = \begin{bmatrix} \frac{1}{s_x} & 0 \\ 0 & \frac{1}{s_y} \end{bmatrix}$  (Pearson correlation then  $\vec{m}_b = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$   $C_b = \begin{bmatrix} 1 \\ \frac{s_{xy}}{s_x s_y} \end{bmatrix}$ 

[on board]

#### Correlation (r) captures dependency

















... but not slope!

















#### Regression (revisited)

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

Optimal regression line slope:

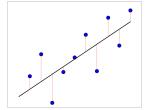
$$\beta = \frac{\vec{x}^T \vec{y}}{\vec{x}^T \vec{x}} = \frac{s_{xy}}{s_x^2}$$

Error variance:

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

Expressed as a proportion of  $\sigma_v^2$ :

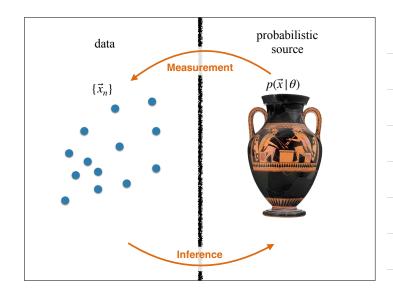
$$\frac{s_e^2}{s_y^2} = 1 - \frac{s_{xy}^2}{s_x^2 s_y^2} = 1 - \boxed{r^2}$$



proportion of data variance explained

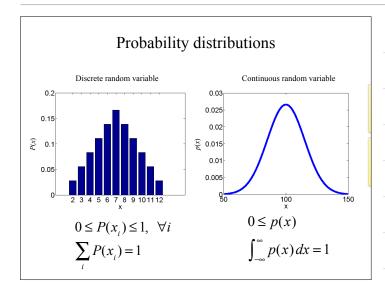
Probability: an abstract mathematical framework for describing random quantities.

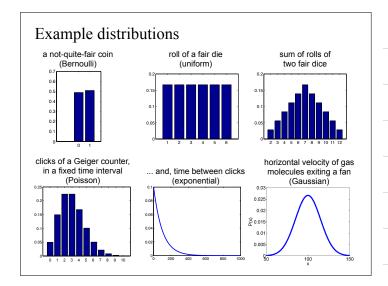
Statistics: use of probability to summarize, analyze, and interpret data. Fundamental to all experimental science.

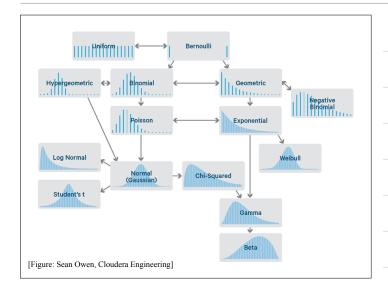


## Univariate Probability (outline)

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



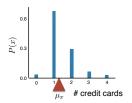




# Expected value (for a discrete random variable)

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

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 values of R.V.)

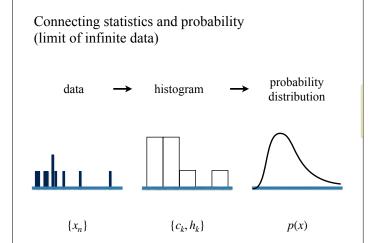
Sample average: an estimate of the expected value:

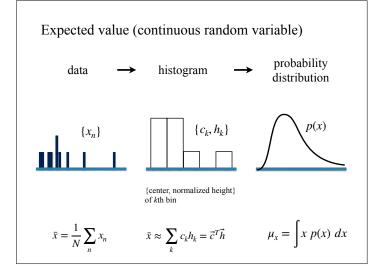
$$ar{f}(x) = rac{1}{N} \sum_{n=1}^N f(x_n)$$
 (sum over data sample

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

#### 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  $(\mu)$  and variance  $(\sigma^2)$ .





#### Expected value (continuous)

$$\mathbb{E}(x) = \int x \ p(x) \ dx$$

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

["mean",  $\mu$ ]

$$\mathbb{E}\left((x-\mu)^2\right) = \int (x-\mu)^2 \ p(x) \ dx \qquad \qquad \text{["variance", $\sigma^2$]}$$

$$= \int x^2 p(x) \ dx - \mu^2 \qquad \qquad \text{[$m_2$ minus $\mu^2$]}$$

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

Note: expectation is an integral, and thus linear, so:

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

#### Transformations of scalar random variables

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

Analogous to sample mean/covariance:

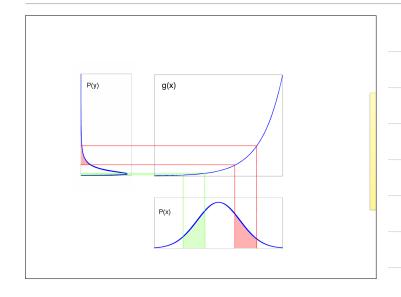
$$\mu_Y = \mathbb{E}(Y) = a\mathbb{E}(X) + b = a\mu_X + b$$

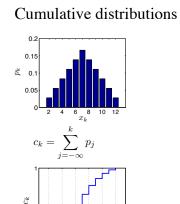
$$\sigma_Y^2 = \mathbb{E}\left(\left(Y - \mu_Y\right)^2\right) = \mathbb{E}\left(\left(aX - a\mu_X\right)^2\right) = a^2\sigma_X^2$$

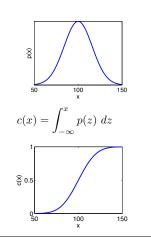
Full distribution: 
$$p_Y(y) = \frac{1}{a} p_X\left(\frac{y-b}{a}\right)$$

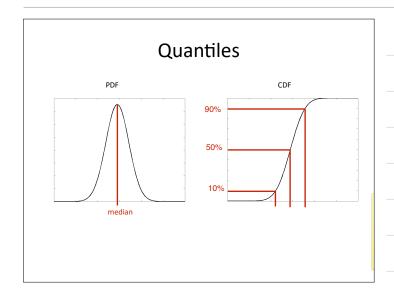
$$Y = g(X)$$
 (assume g is "monotonic" - i.e., derivative > 0)

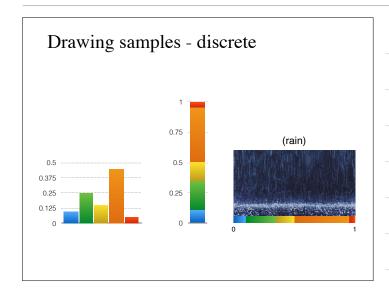
$$p_{Y}(y) = \frac{p_{X}(g^{-1}(y))}{g'(g^{-1}(y))}$$











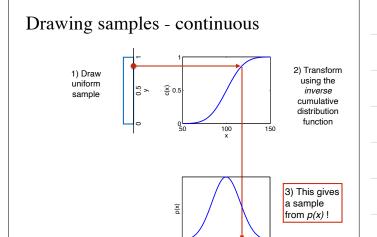
# Drawing samples - continuous 3) Result is uniformly distributed! 2) Transform using the cumulative distribution function

$$p_{y}(y) = \frac{p_{x}(c^{-1}(y))}{c'(c^{-1}(y))}$$

$$= \frac{p_{x}(c^{-1}(y))}{p_{x}(c^{-1}(y))}$$

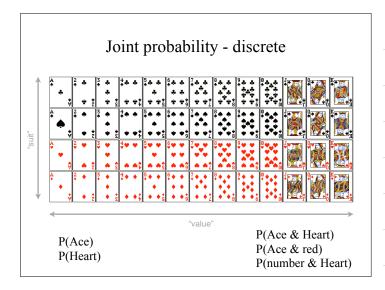
$$= 1, \quad 0 \le y \le 1$$

1) Sample from *p(x)* 

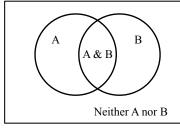


# Multi-variate probability (outline)

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

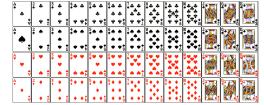


# Conditional probability



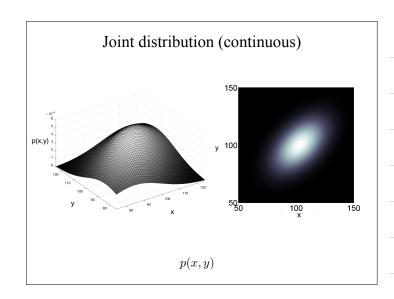
p(A | B) = probability of A given that B is asserted to be true =  $\frac{p(A \& B)}{p(B)}$ 

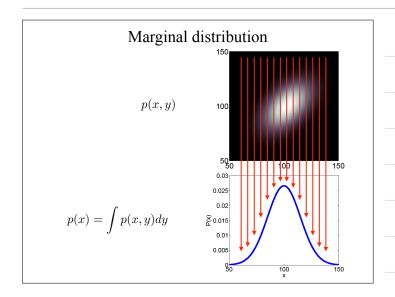
# Conditional probability - discrete

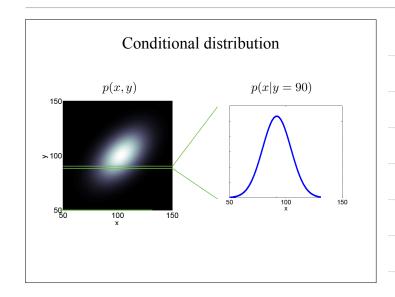


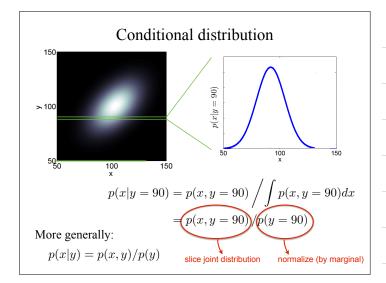


P(Ace | Heart)
P(Ace | red)
P(number | Heart)









# 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.

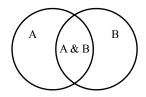
Dear Sir,

Read Dec. 25, Town fend you an effay which I have 1765. I found among the papers of our decased 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



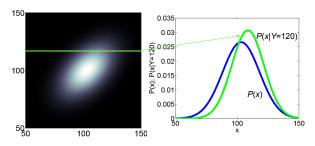
p(A|B) = probability of A given that B is asserted to be true =  $\frac{p(A \& B)}{a}$ 

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

$$= p(A)p(B \mid A)$$

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

#### Conditional vs. marginal

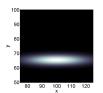


In general, the marginals for different Y values differ. When are they they same? In particular, when are all conditionals equal to the marginal?

## 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. 
$$\vec{x}$$
,  $\overrightarrow{\mu}_x = \mathbb{E}(\vec{x})$ ,  $C_x = \mathbb{E}\left((\vec{x} - \overrightarrow{\mu}_x)(\vec{x} - \overrightarrow{\mu}_x)^T\right)$ 

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

analogous to results for sample mean/covariance:

$$\overrightarrow{\mu}_{y} = \mathbb{E}(M(\overrightarrow{x} - \overrightarrow{a}))$$

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

$$= M(\overrightarrow{\mu}_{x} - \overrightarrow{a})$$

$$C_{y} = \mathbb{E}((M(\overrightarrow{x} - \overrightarrow{\mu}_{x}))(M(\overrightarrow{x} - \overrightarrow{\mu}_{x}))^{T})$$

$$= M\mathbb{E}((\overrightarrow{x} - \overrightarrow{\mu}_{x}))(\overrightarrow{x} - \overrightarrow{\mu}_{x}))^{T})M^{T}$$

$$= MC_{x}M^{T}$$

#### 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:

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

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

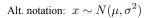
$$p_Z(z)$$
 is the *convolution* of  $p_X(x)$  and  $p_Y(y)$ 

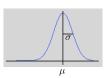
[on board]

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

One-dimensional:

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





Multi-dimensional:

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

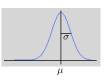


mean: [0.2, 0.8] cov: [1.0 -0.3; -0.3 0.4]

#### Gaussian properties

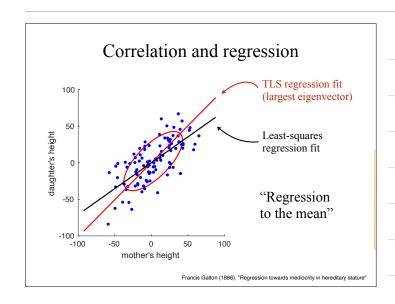
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

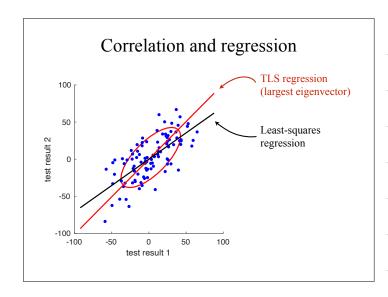
$$p(\vec{x}) = \frac{1}{\sqrt{(2\pi)^N |C|}} \; e^{-(\vec{x} - \vec{\mu})^T C^{-1} (\vec{x} - \vec{\mu})/2}$$

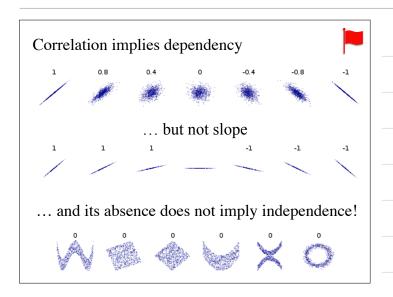


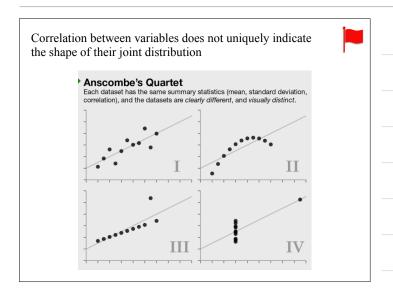
- joint density of indep Gaussian RVs is elliptical [easy]
- conditionals of a Gaussian are Gaussian [easy]
- marginals of a Gaussian are Gaussian [easy]
- product of two Gaussian dists is Gaussian [easy]
- sum of independent Gaussian RVs is Gaussian [moderate]
- the most random (max entropy) density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

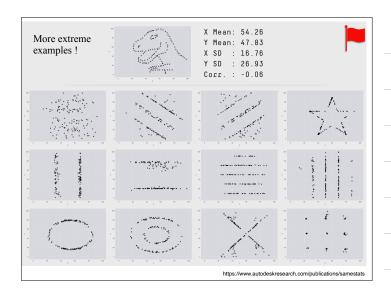
Generalized marginals of a Gaussian 
$$\vec{x} \sim N(\vec{\mu}_x, C_x)$$
 
$$z = \hat{u}^T \vec{x}$$
 
$$p(z) \text{ is Gaussian, with:}$$
 
$$\mu_z = \hat{u}^T \vec{\mu}_x$$
 
$$\sigma_z^2 = \hat{u}^T C_x \hat{u}$$

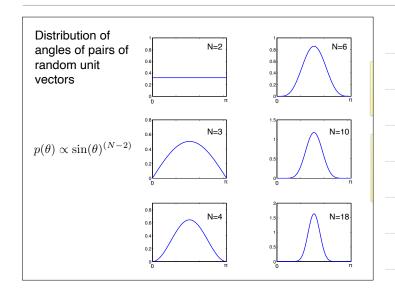


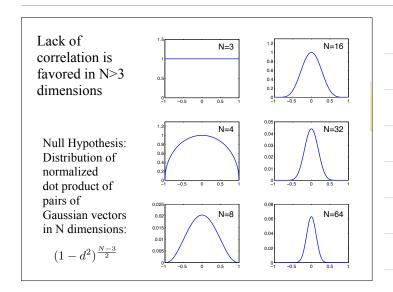


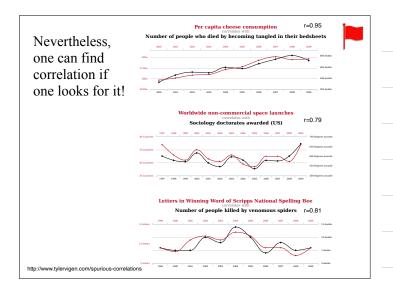










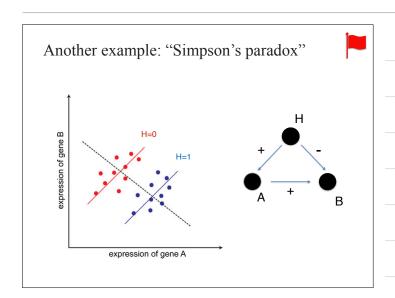


#### Covariation/correlation does not imply causation

- 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: 
$$p(A, B \mid H) = p(A \mid H)p(B \mid H)$$

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



# Milton Friedman's Thermostat O = outside temperature (assumed cold) I = inside temperature (ideally, constant) E = energy used for heating Observed statistics, P=C-1: Statistical observations: O and I uncorrelated I and E uncorrelated E

O and E anti-correlated

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 does *not* imply data lie near a line/plane/hyperplane.
- Independent implies uncorrelated. But uncorrelated does *not* imply independent.
- Correlation does *not* imply causation (and often arises from hidden common factors).
- Correlation is a **descriptive statistic**, and does not eliminate the need for reasoning/experiments/models!