



GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

Deep Learning

Lecture 2: Math basics

Alexander Ecker
Institut für Informatik, Uni Göttingen



<https://alexanderecker.wordpress.com>

– Credit: slides largely based on Ian Goodfellow's and Chris Bishop's slides –

Recommended literature

Primary book:

- **Goodfellow, Bengio, Courville: Deep Learning**
<https://www.deeplearningbook.org>

Good general machine learning book:

- Christopher Bishop: Pattern Recognition and Machine Learning
<https://www.microsoft.com/en-us/research/people/cmbishop/prml-book/>

Good resource for Deep Reinforcement Learning (last week only):

- OpenAI: Spinning Up in Deep RL
<https://openai.com/blog/spinning-up-in-deep-rl/>

Homework assignments

First homework assignment on Thursday

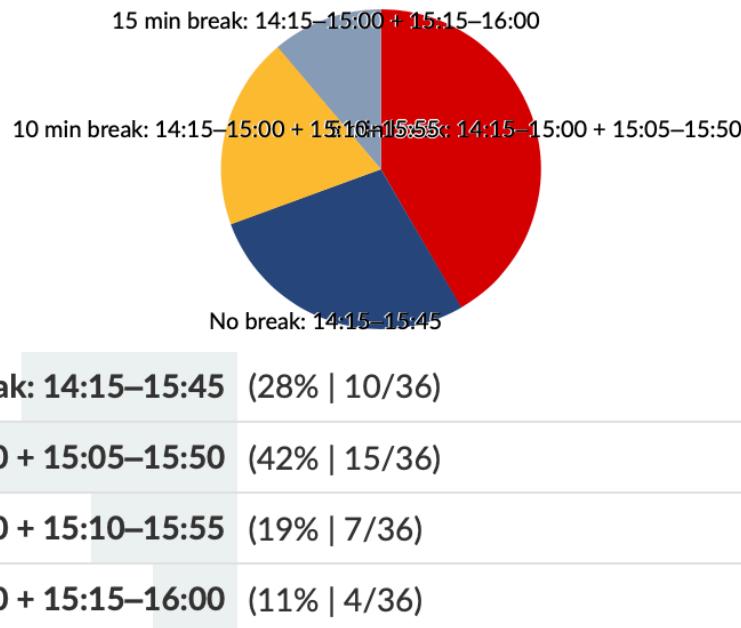
Form teams of 2–3 students

- No individuals please
- You'll present your solution to the tutors in tutorial (5 min)
- It's ok to split work, but everybody needs to be able to explain everything

Details on submission on Thursday

To break or not to break?

📊 Should we take a break after 45 minutes or do 90 min straight?



$$\mathbb{E}[\text{break length}] = 5 \text{ min } 39 \text{ sec} \rightarrow \mathbf{5 \text{ min break}}$$

Topics today

Linear algebra

Probability theory

Information theory

Linear algebra

Notation

Scalar: a single number

$$a, n, x \in \mathbb{R}$$

0d tensor 

Vector: 1d array of numbers

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$$

1d tensor 

Matrix: 2d array of numbers

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

2d tensor 

Tensor: N -dim array of numbers $\mathbf{W} = (W_{i_1 i_2 \dots i_N})$

3d tensor 

Matrix transpose

$$(\mathbf{A}^T)_{ij} = A_{ji}$$

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \end{bmatrix} \Rightarrow \mathbf{A}^T = \begin{bmatrix} A_{11} & A_{21} \\ A_{12} & A_{22} \\ A_{13} & A_{23} \end{bmatrix}$$

The matrix transpose is a mirror image across the diagonal

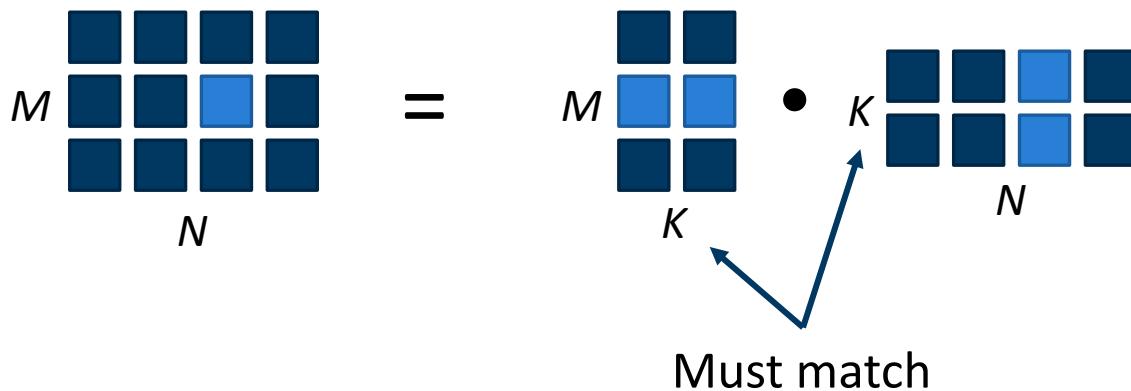
$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

Dot product

(a.k.a. inner product, matrix product)

$$C = AB$$

$$C_{ij} = \sum_{k=1}^K A_{ik}B_{kj}$$



Identity matrix

$$I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$Ix = x$$

$$IA = A = AI$$

Linear system of equations

$$A\boldsymbol{x} = \boldsymbol{b}$$

expands to

$$A_{1:}\boldsymbol{x} = b_1$$

$$A_{2:}\boldsymbol{x} = b_2$$

...

$$A_{m:}\boldsymbol{x} = b_m$$

$$A_{2:} \quad A_{:,1}$$


The image shows two 3x2 grids of colored squares. The left grid, labeled $A_{2:}$, has three rows and two columns. The top row consists of dark blue squares, the middle row of light blue squares, and the bottom row of dark blue squares. The right grid, labeled $A_{:,1}$, also has three rows and two columns. The first column (leftmost) consists of light blue squares, and the second column (rightmost) consists of dark blue squares.

Solving linear systems of equations

A linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ can have...

- No (exact) solution

\mathbf{A} is *tall*

A diagram illustrating a tall matrix \mathbf{A} . It consists of a vertical column of four dark blue squares. To its right is a multiplication sign (\cdot) followed by a horizontal vector \mathbf{x} , which is represented by two dark blue squares side-by-side. To the right of the multiplication sign is an equals sign ($=$) followed by another horizontal vector \mathbf{b} , also represented by two dark blue squares side-by-side.

- Many solutions

\mathbf{A} is *fat*

A diagram illustrating a fat matrix \mathbf{A} . It consists of a horizontal row of four dark blue squares. To its right is a multiplication sign (\cdot) followed by a vertical vector \mathbf{x} , which is represented by two dark blue squares stacked vertically. To the right of the multiplication sign is an equals sign ($=$) followed by another vertical vector \mathbf{b} , also represented by two dark blue squares stacked vertically.

- Exactly one solution

\mathbf{A} is *square and invertible*

A diagram illustrating a square and invertible matrix \mathbf{A} . It consists of a square arrangement of four dark blue squares. To its right is a multiplication sign (\cdot) followed by a vertical vector \mathbf{x} , which is represented by two dark blue squares stacked vertically. To the right of the multiplication sign is an equals sign ($=$) followed by another vertical vector \mathbf{b} , also represented by two dark blue squares stacked vertically.

Matrix inversion

Matrix inverse

$$A^{-1}A = I = AA^{-1}$$

Solving linear system via matrix inverse

$$\begin{aligned} Ax &= b \\ A^{-1}Ax &= A^{-1}b \\ Ix &= A^{-1}b \\ x &= A^{-1}b \end{aligned}$$

Note: don't do it like this in practice!

Numerically unstable, but useful for analytical analysis

Norms

Measures the “length” of a vector

If $\|\cdot\|$ represents a norm, it must hold:

- (1) $\|x\| = 0 \Rightarrow x = 0$
- (2) $\|x + y\| < \|x\| + \|y\|$ (triangle inequality)
- (3) $\|a \cdot x\| = |a| \cdot \|x\|$

L_p norm (most popular: L_2):

$$\|x\|_p = \left(\sum_i |x_i|^p \right)^{1/p}$$

L_1 norm: $\|x\|_1 = \sum |x_i|$

Maximum (L_∞) norm: $\|x\|_\infty = \max_i |x_i|$

Special matrices and vectors

Unit vector:

$$\|x\| = 1$$

Symmetric matrix:

$$A^T = A$$

Orthogonal matrix:

$$\begin{aligned} A^T A &= A A^T = I \\ A^T &= A^{-1} \end{aligned}$$

Eigenvalue decomposition

Eigenvector \boldsymbol{v} and eigenvalue λ :

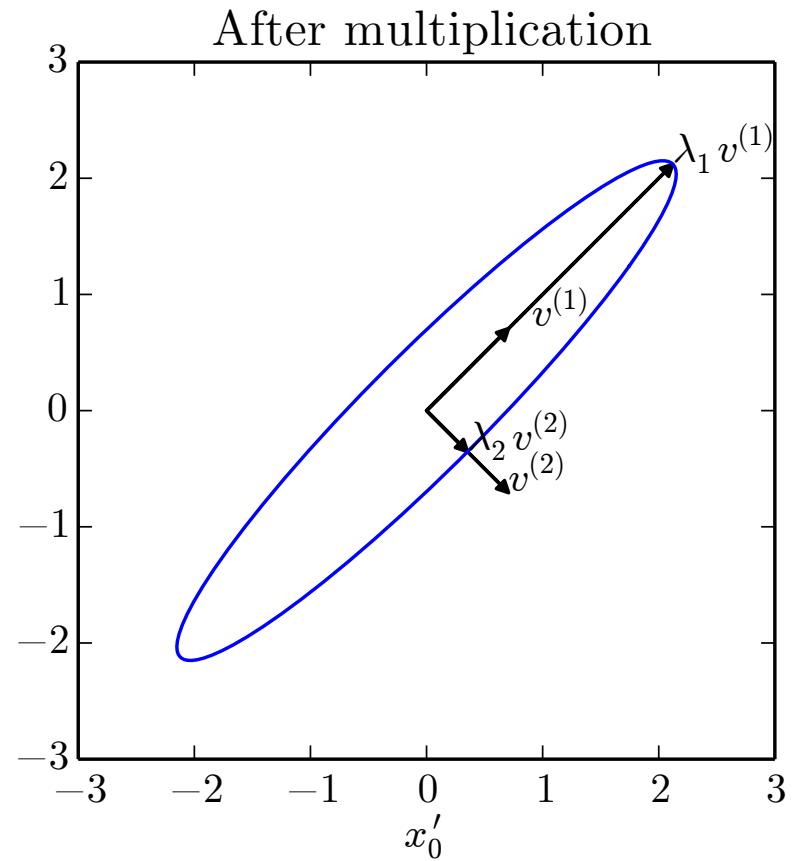
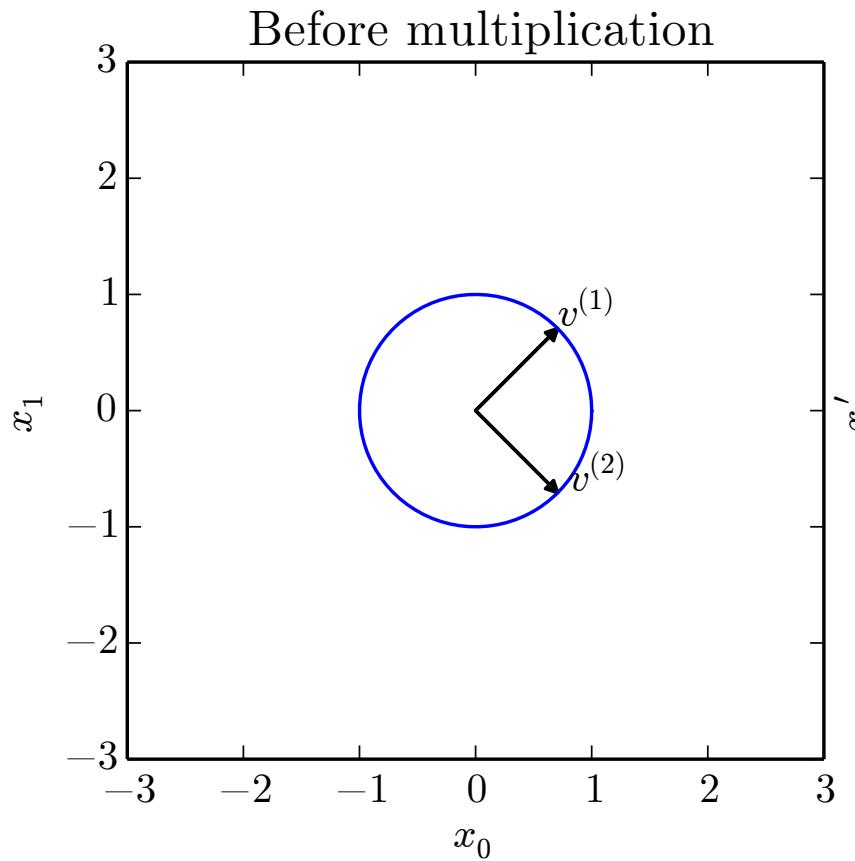
$$A\boldsymbol{v} = \lambda\boldsymbol{v}$$

Eigenvalue decomposition of a diagonalizable matrix A

$$\begin{aligned} A &= V\Lambda V^T \\ V^T V &= I \quad \Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} \end{aligned}$$

Every real symmetric matrix is diagonalizable

Eigenvalue decomposition graphically



Singular value decomposition

Similar to eigenvalue decomposition

More general: matrix need not be square

$$A = USV^T$$

$$U^T U = I \quad V^T V = I \quad S = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{bmatrix}$$

Singular values

$$\begin{matrix} A \\ \begin{matrix} \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \end{matrix} \end{matrix} = \begin{matrix} U \\ \begin{matrix} \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \end{matrix} \end{matrix} \cdot \bullet \cdot \begin{matrix} S \\ \begin{matrix} \blacksquare & \square & \blacksquare \\ \square & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{matrix} \end{matrix} \cdot \bullet \cdot \begin{matrix} V^T \\ \begin{matrix} \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \end{matrix} \end{matrix}$$

Matrix pseudoinverse

(a.k.a. Moore-Penrose inverse)

$$AA^\dagger A = A \quad A^\dagger AA^\dagger = A^\dagger$$

Symmetric

$$(A^\dagger A)^T = A^\dagger A \quad (AA^\dagger)^T = AA^\dagger$$

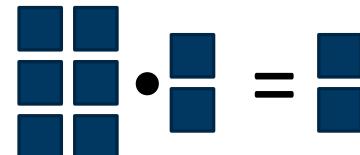
For linearly independent columns of A :

$$A^\dagger = (A^T A)^{-1} A^T$$
$$A^\dagger A = I$$



Least-squares solution to overdetermined
system of linear equations:

$$x = A^\dagger b$$



Matrix pseudoinverse

System of linear equations:

$$\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}$$

1. System has exactly one solution:

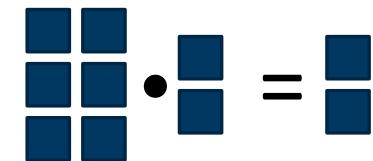
same as $\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$

2. System has no solution:

\mathbf{x} with the smallest quadratic error $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$

3. System has many solutions:

\mathbf{x} with the smallest L_2 norm $\|\mathbf{x}\|_2$



A diagram illustrating a system of linear equations. It shows a 3x3 grid of dark blue squares representing matrix \mathbf{A} , followed by a dot, a 2x2 grid of dark blue squares representing vector \mathbf{b} , followed by an equals sign, and a 3x2 grid of dark blue squares representing the solution vector \mathbf{x} . The result is a square 3x2 matrix.



A diagram illustrating a system of linear equations. It shows a 3x3 grid of dark blue squares representing matrix \mathbf{A} , followed by a dot, a 2x2 grid of dark blue squares representing vector \mathbf{b} , followed by an equals sign, and a 3x2 grid of dark blue squares representing the solution vector \mathbf{x} . The result is a rectangular 3x2 matrix.

Trace

Sum of the diagonal elements of a matrix

$$\text{Tr}(\mathbf{A}) = \sum_i A_{ii}$$

Trace permutes cyclically

$$\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{CAB}) = \text{Tr}(\mathbf{BCA})$$

Determinant

For square matrices:

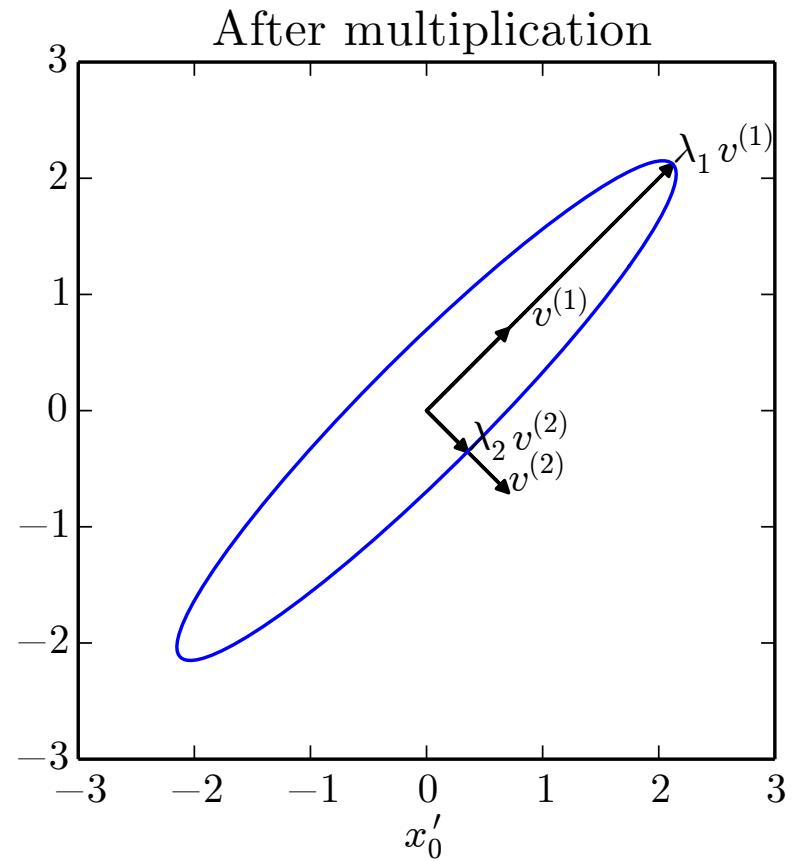
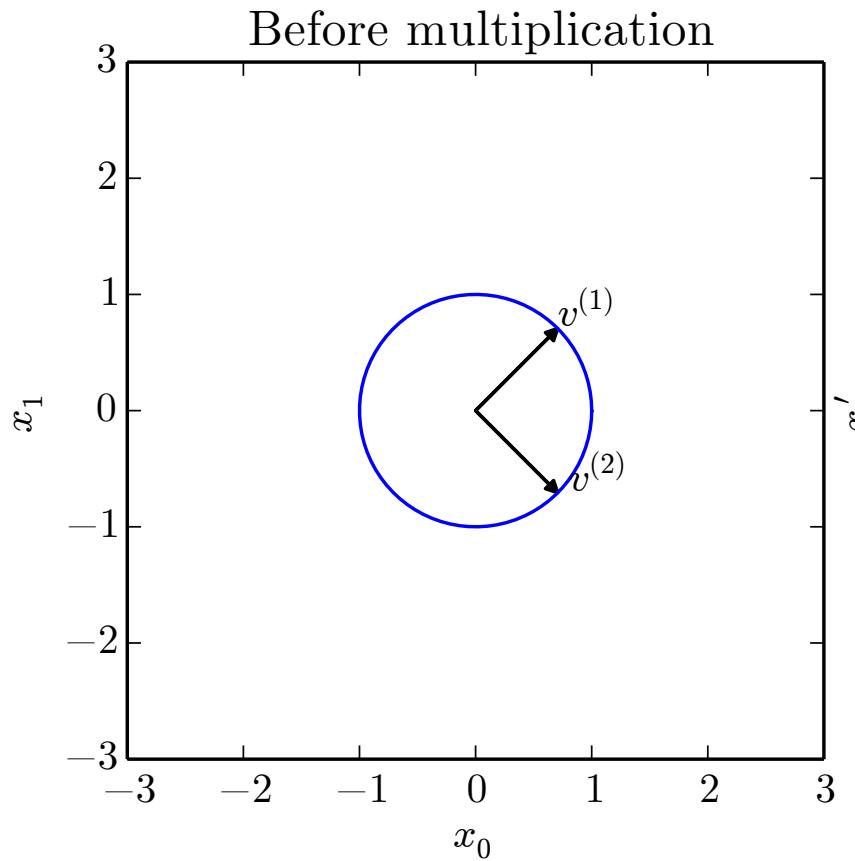
$$\det(\mathbf{A}) = \prod_k \lambda_k \quad \xleftarrow{\text{Eigenvalues}}$$

Absolute value of the determinant measures how much multiplication by a matrix expands or contracts the space

Matrix \mathbf{A} is invertible if $\det(\mathbf{A}) \neq 0$

Volume-conserving transformations have $\det(\mathbf{A}) = 1$

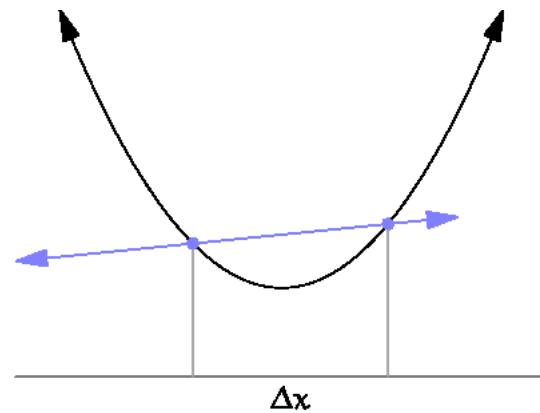
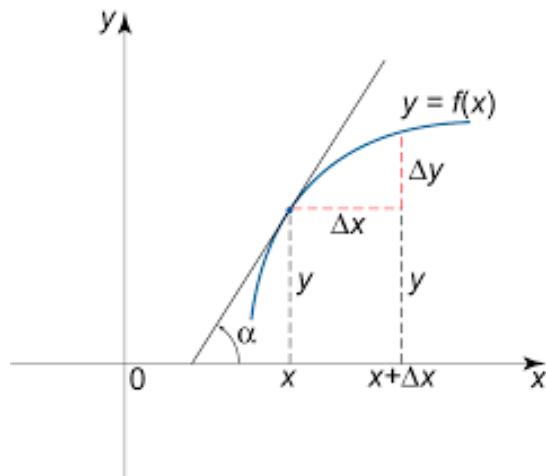
Eigenvalue decomposition graphically



Derivative

The derivative $f'(x)$ of a function $f(x)$ is its slope at x

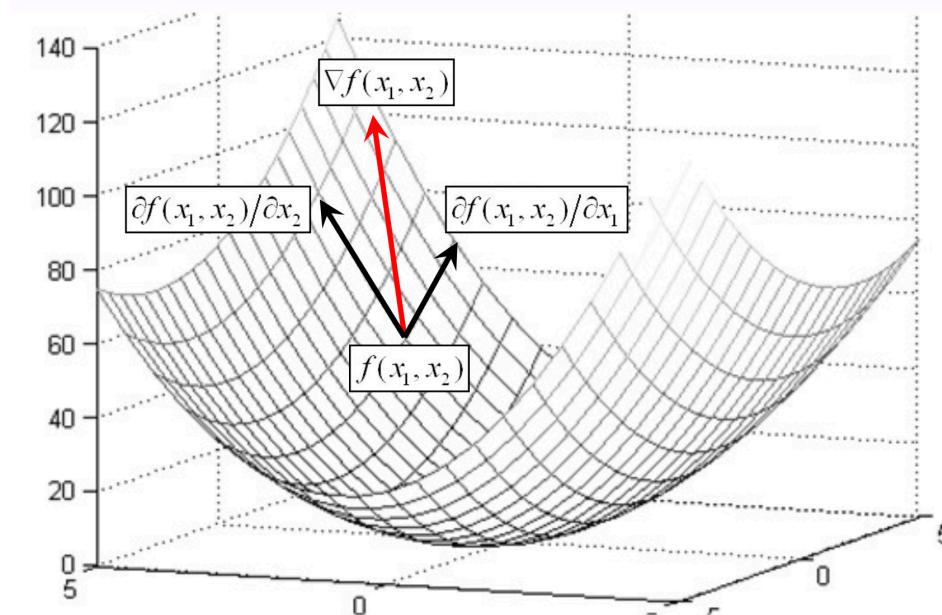
$$f'(x) = \frac{\partial f}{\partial x}(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$



Gradient

“Derivative” of a function $f(\mathbf{x})$ with vector-valued inputs \mathbf{x}
Points into direction of steepest ascent

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}) = [f'(x_1), \dots, f'(x_n)]$$



Jacobian

Gradient matrix of function $f(\mathbf{x})$ with vector-valued input and output

$$J_{\mathbf{x}} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}) \\ \vdots \\ \nabla_{\mathbf{x}} f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{x_1} & \dots & \frac{\partial f_1}{x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{x_1} & \dots & \frac{\partial f_m}{x_n} \end{bmatrix}$$

Chain rule

Derivate of compositions

$$f'(g(x)) = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}(x)$$

Example $f(x) = \sin(x^2)$

$$f'(x) = \frac{\partial \sin(g)}{\partial g} \cdot \frac{\partial x^2}{\partial x} = \cos(x^2) \cdot 2x$$

Vector-valued functions:

$$\nabla_x f(\mathbf{g}(x)) = \frac{\partial f}{\partial \mathbf{g}} \cdot \frac{\partial \mathbf{g}}{\partial x}(x)$$

$$\begin{matrix} \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \end{matrix} = \begin{matrix} \blacksquare & \blacksquare & \blacksquare \end{matrix} \bullet \begin{matrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{matrix} \quad \begin{matrix} \blacksquare \\ \blacksquare \end{matrix}$$

$\nabla_x \qquad \nabla_g \qquad J_x \qquad x$

Linear algebra cheat sheet

The Matrix Cookbook

http://www2.imm.dtu.dk/pubdb/views/publication_details.php?id=3274

Contains lots of useful linear algebra identities

Probability theory

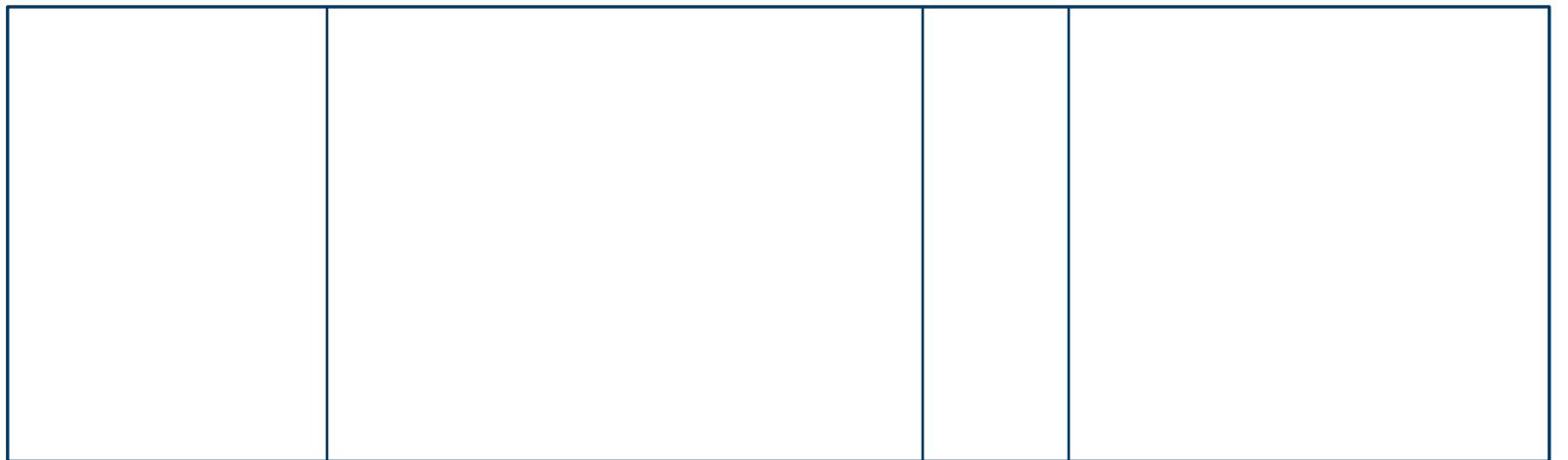
Why probability theory?

Reasoning about the future state of the world requires dealing with uncertainty

- Inherent stochasticity in the system
- Observation noise
- Incomplete observability
- Incomplete modeling

Events

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



0

0.2

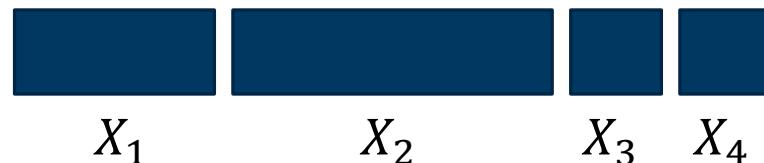
0.6

0.7

1

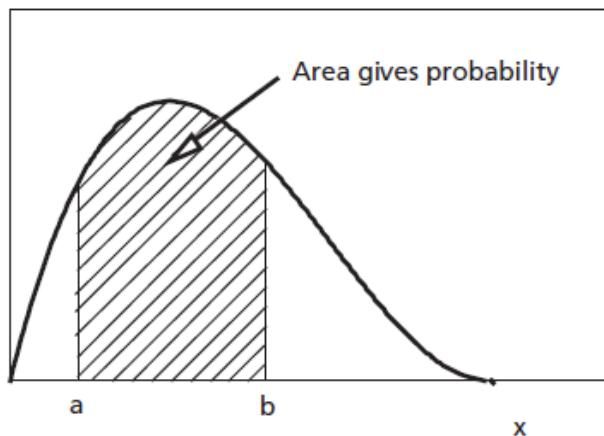
Discrete vs. continuous random variables

Discrete



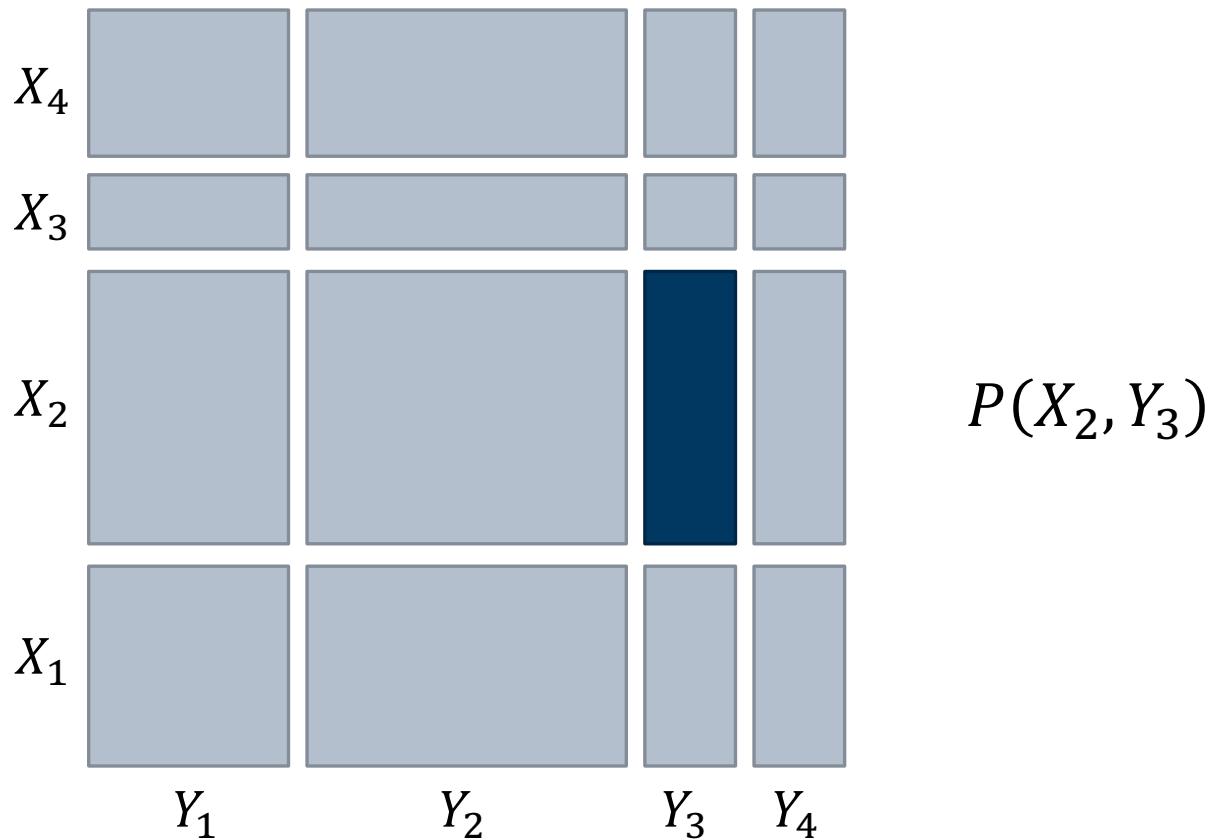
Probability mass function
 $P(X = X_1)$

Continuous

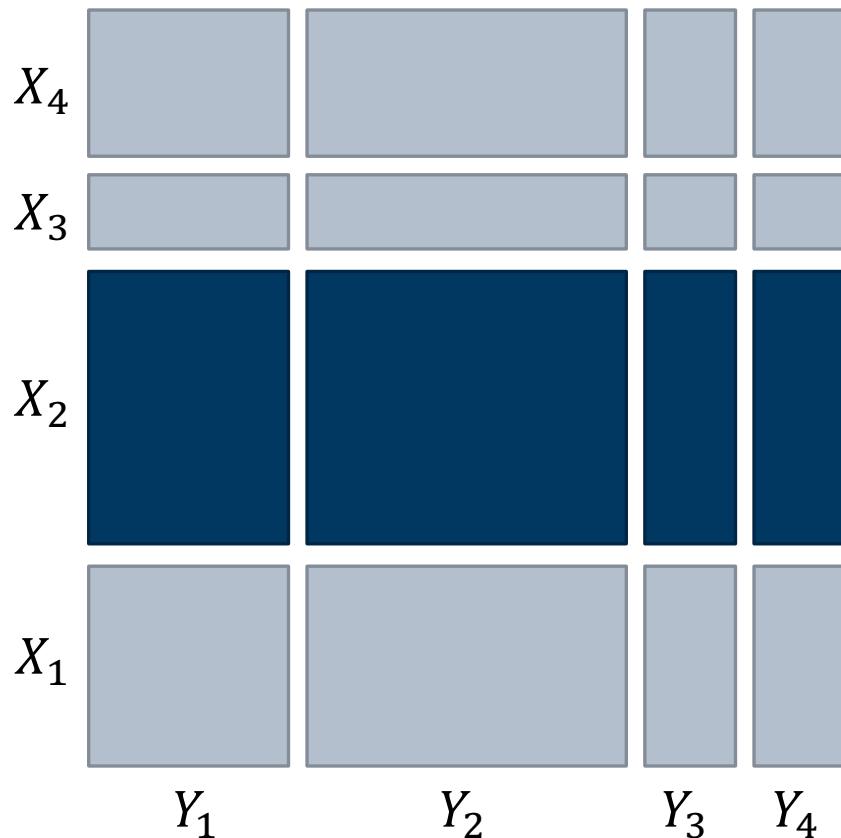


Probability density function $p(x)$
 $P(X \in [a, b]) = \int_a^b p(x)dx$

Joint probability

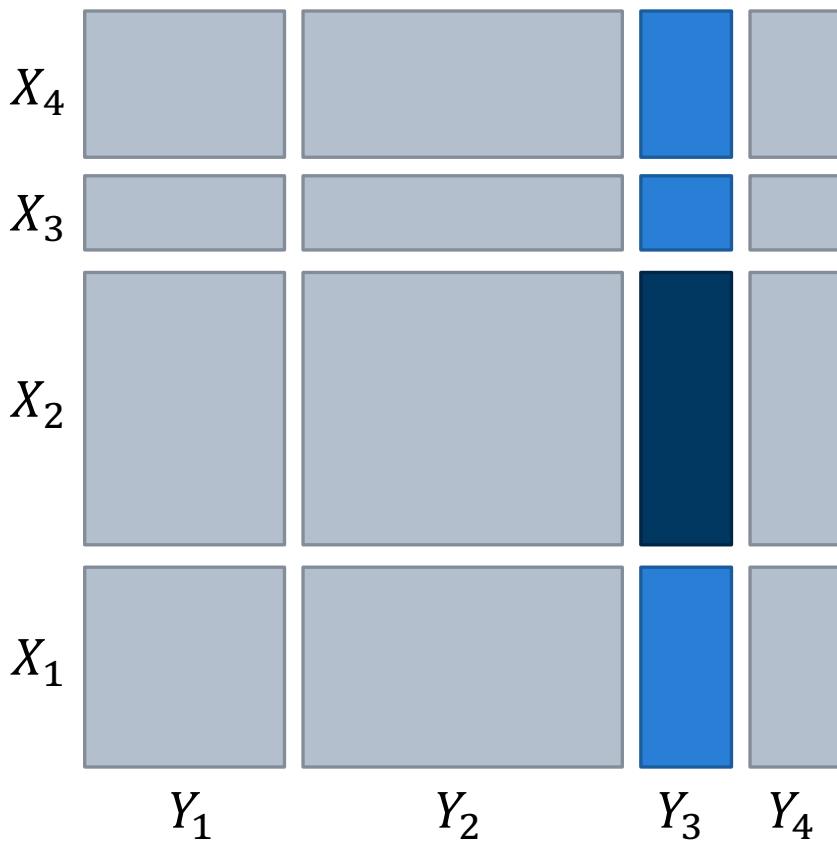


Marginal probability



$$P(X_2) = \sum_k P(X_3, Y_k)$$

Conditional probability



$$P(X_2|Y_3) = \frac{P(X_2, Y_3)}{P(Y_3)}$$

Bayes' rule

$$\Rightarrow \begin{aligned} P(X, Y) &= P(X|Y)P(Y) \\ P(Y, X) &= P(Y|X)P(X) \end{aligned}$$

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Posterior \propto Likelihood \times Prior

Example: Bayes' rule

Morse code (long “–” or short “•”) transmitted through a noisy line

$$P(S = “–”) = \frac{4}{7} \quad P(S = “•”) = \frac{3}{7} \quad P(\text{flip}) = \frac{1}{8}$$

A “•” was received. What is the probability that a “•” was sent?

$$\begin{aligned} P(R = “•”) &= P(R = “•”|S = “•”)P(S = “•”) + \\ &\quad P(R = “•”|S = “–”)P(S = “–”) \\ &= \frac{7}{8} \cdot \frac{3}{7} + \frac{1}{8} \cdot \frac{4}{7} = \frac{25}{56} \end{aligned}$$

$$P(S = “•”|R = “•”) = \frac{P(R = “•”|S = “•”)P(S = “•”)}{P(R = “•”)} = \frac{\frac{7}{8} \cdot \frac{3}{7}}{\frac{25}{56}} = \frac{21}{25} = 84\%$$

Independence

X and Y are *independent* if and only if

$$P(X, Y) = P(X)P(Y) \quad \text{for all } X, Y$$

X and Y are *conditionally independent* if and only if

$$P(X, Y|Z) = P(X|Z)P(Y|Z) \quad \text{for all } X, Y, Z$$

Expectation

Discrete random variables:

$$\mathbb{E}[X] = \sum_k P(X = x_k)x_k$$

Continuous random variables

$$\mathbb{E}[X] = \int_{\mathcal{X}} p(x)x \, dx$$

Expectation is a linear operation:

$$\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$$

Variance and covariance

Variance measures the spread of a distribution
(average squared deviation from the mean)

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

Standard deviation: $\text{SD}[X] = \sqrt{\text{Var}[X]}$

Covariance measures the co-fluctuations of two random variables

$$\begin{aligned}\text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

Covariance matrix $\mathbf{C}_{ij} = \text{Cov}[X_i, X_j]$

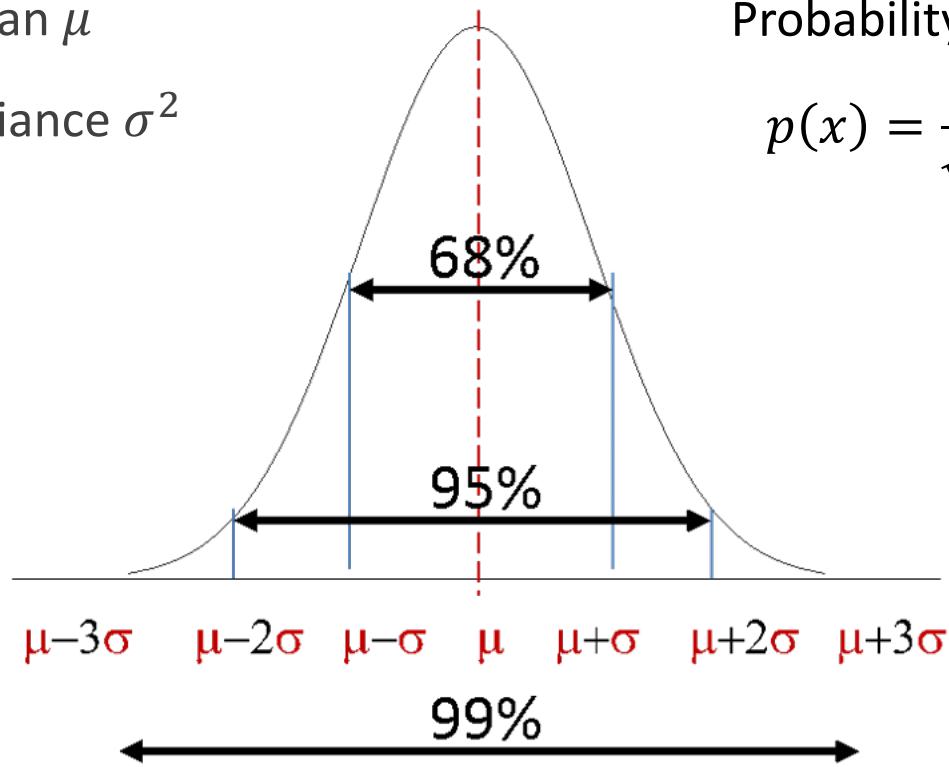
Gaussian distribution

Mean μ

Variance σ^2

Probability density function

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



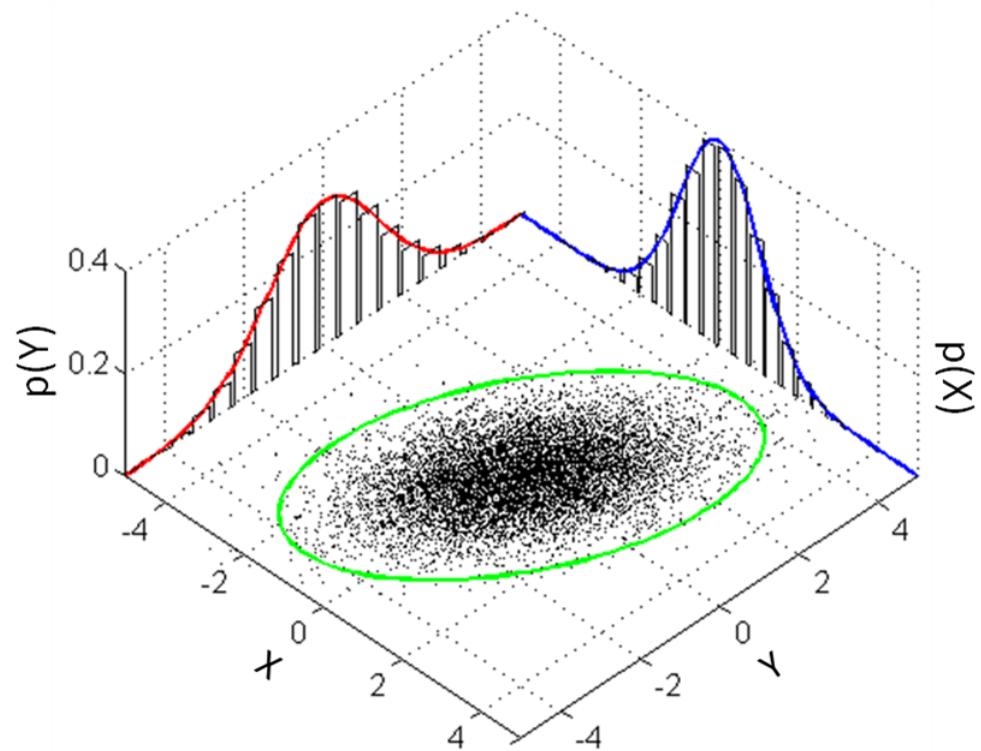
Multivariate normal distribution

Probability density function

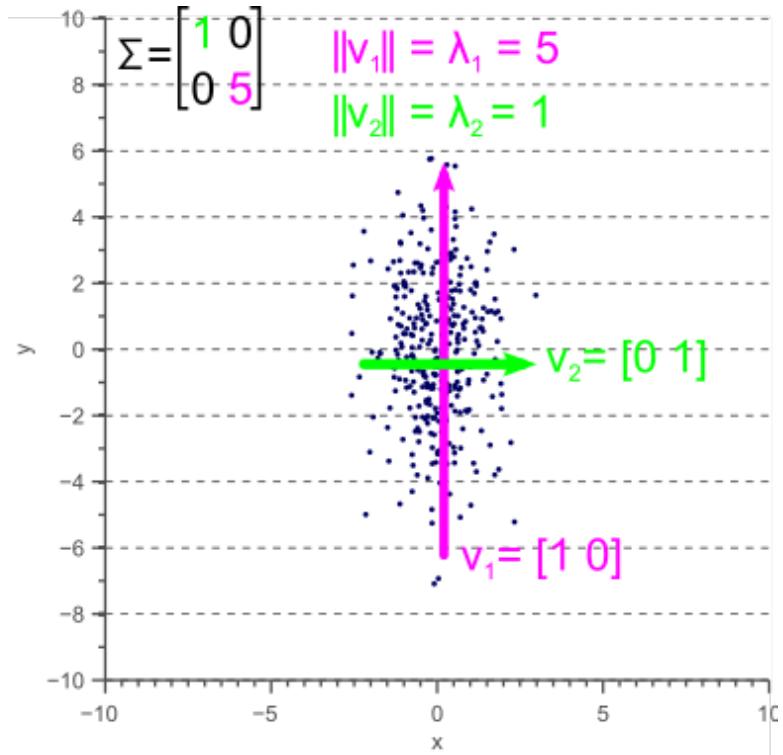
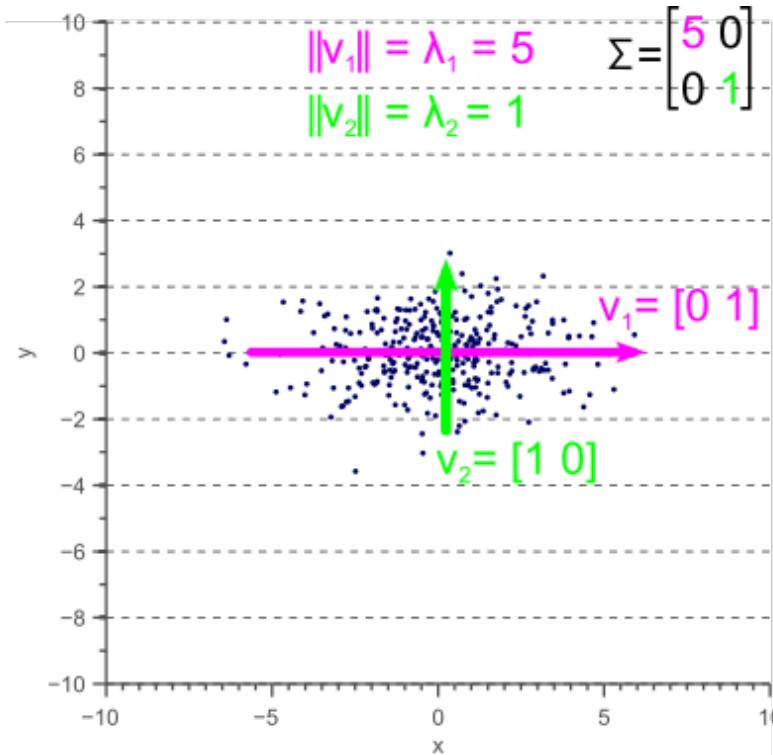
$$p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n |\mathbf{C}|}} \exp(-(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}))$$

Mean $\boldsymbol{\mu}$

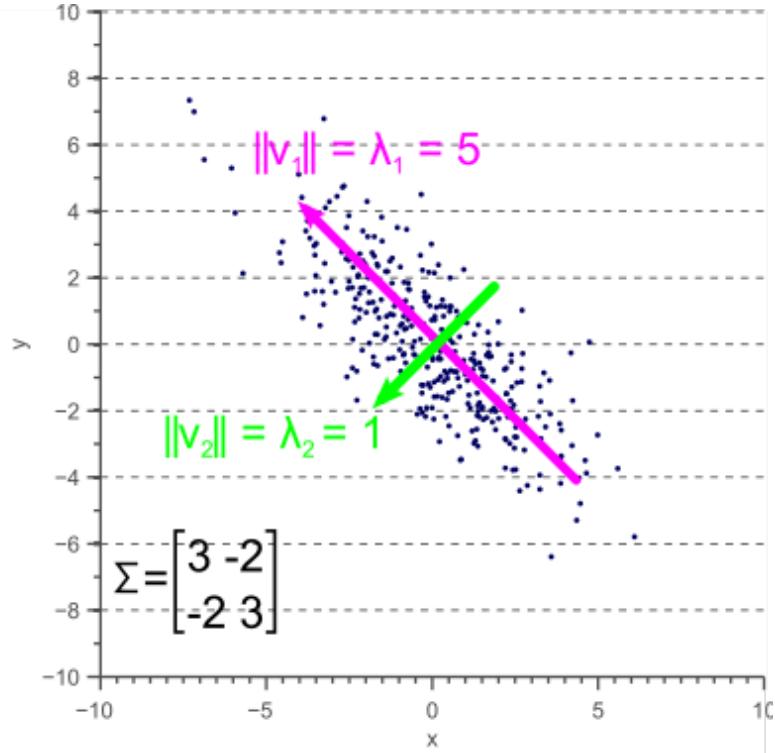
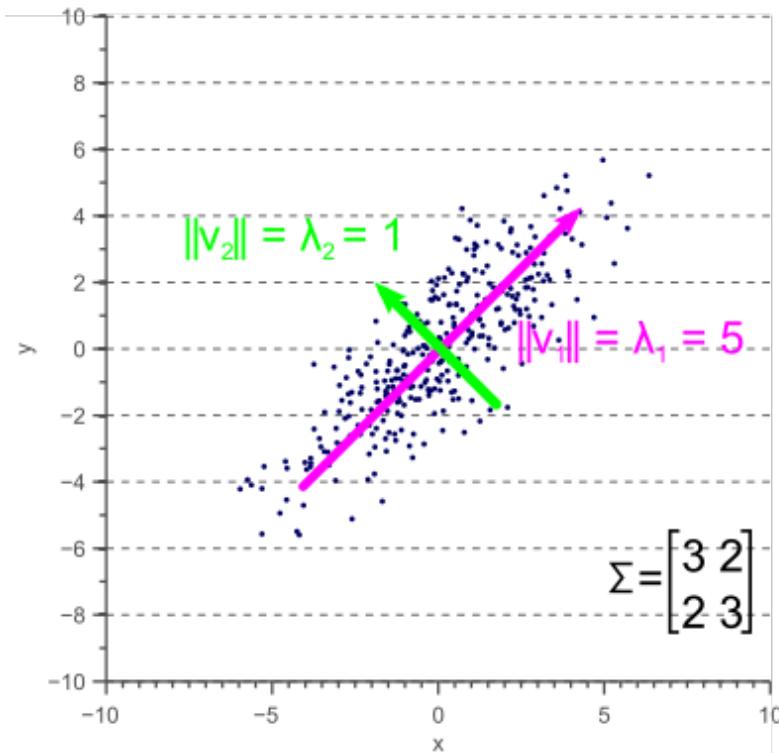
Covariance matrix \mathbf{C}



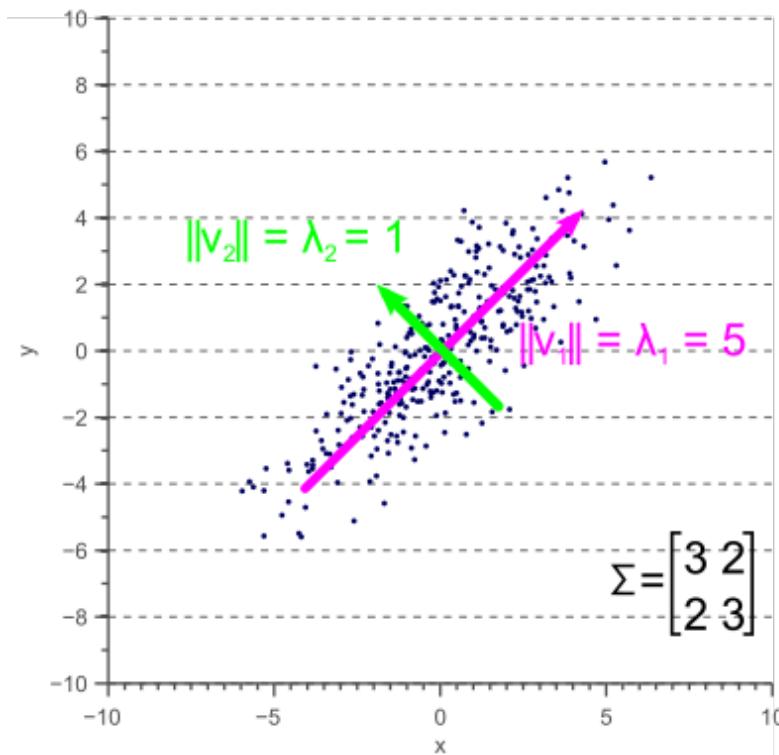
Uncorrelated 2d Gaussian



Correlated 2d Gaussian



Correlated 2d Gaussian



Eigenvalue decomposition of covariance matrix \mathbf{C}

$$\mathbf{C} = \mathbf{V}\Lambda\mathbf{V}^T$$

$$\mathbf{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix}$$

→ Principal Component Analysis (PCA)

Bernoulli + multinomial distribution

Bernoulli distribution models binary random variable

$$\begin{aligned}P(X = 1) &= p \\P(X = 0) &= 1 - p\end{aligned}$$

$$P(X = x) = p^x(1 - p)^{1-x}$$

Multinomial distribution ($N = 1$) models categorical random variable

$$P(X = \mathbf{x}) = \prod_k p_k^{x_k} \quad \text{where } \mathbf{p} \in [0,1]^K \quad \text{and } p_K = 1 - \sum_{k=1}^K p_k$$

Often used in classification problems

Information theory

Entropy

$$H[x] = - \sum_x p(x) \log_2 p(x)$$

Important quantity in

- coding theory
- statistical physics
- machine learning

Entropy

Coding theory: discrete x with 8 possible states; how many bits to transmit the state of x ?

All states equally likely

$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$

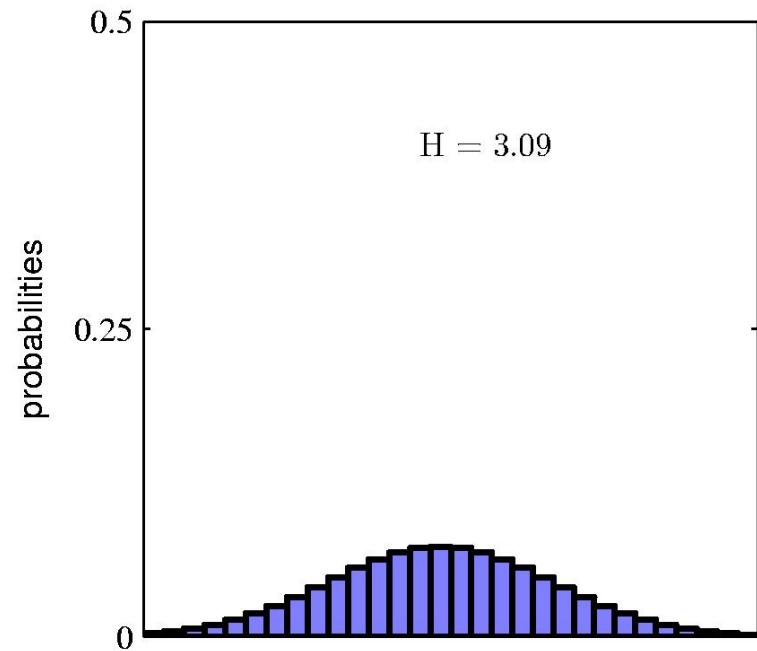
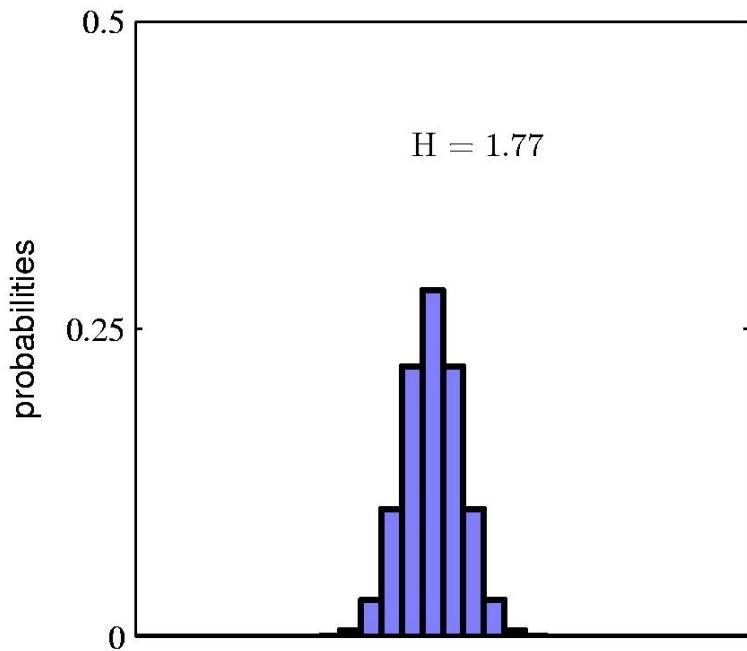
Entropy

x	a	b	c	d	e	f	g	h
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$
code	0	10	110	1110	111100	111101	111110	111111

$$\begin{aligned} H[x] &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} \\ &= 2 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{average code length} &= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 \\ &= 2 \text{ bits} \end{aligned}$$

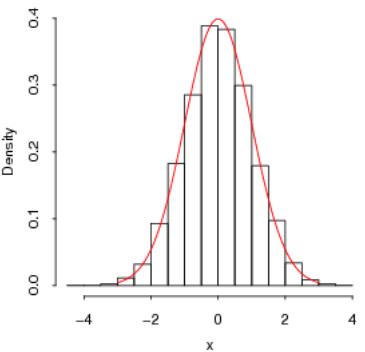
Entropy



Differential Entropy

Put bins of width Δ along the real line

$$\lim_{\Delta \rightarrow 0} \left\{ - \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x)$$



Differential entropy maximized (for fixed σ^2) when

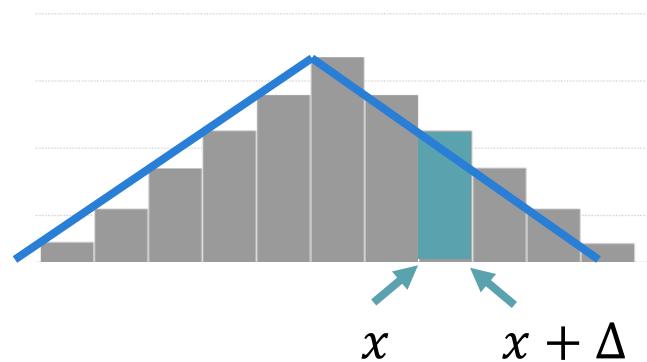
$$p(x) = \mathcal{N}(x|\mu, \sigma^2)$$

in which case $H[x] = \frac{1}{2} \{1 + \ln(2\pi\sigma^2)\}$.

Differential Entropy

Put bins of width Δ along the real line

$$\begin{aligned} H_\Delta &= - \sum_i p(x_i) \Delta \ln(p(x_i) \Delta) \\ &= - \sum_i p(x_i) \Delta \ln p(x_i) - \ln \Delta \end{aligned}$$



$$H[x] = \lim_{\Delta \rightarrow 0} \left\{ - \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x) dx$$

Conditional Entropy

Suppose we have a joint distribution $p(x, y)$ from which we draw pairs of values of x and y

If a value of x is already known, then the additional information needed to specify the corresponding value of y is given by $-\ln p(y|x)$.

$$H[y|x] = - \iint p(y, x) \ln p(y|x) dy dx$$

$$H[x, y] = H[y|x] + H[x]$$

The Kullback-Leibler Divergence

Suppose we have modeled some unknown distribution $p(x)$ using an approximating distribution $q(x)$

The KL divergence quantifies additional amount of information required to specify the value of x as a result of using $q(x)$ instead of $p(x)$

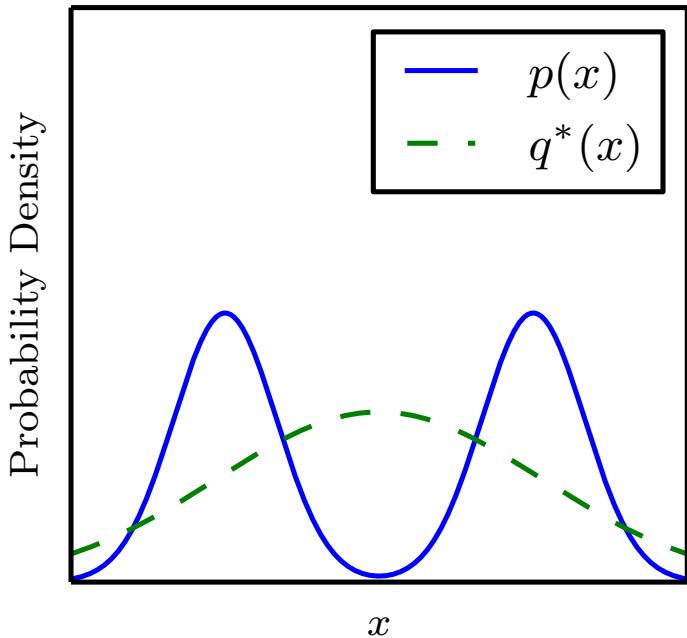
$$\begin{aligned} \text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \end{aligned}$$

$$\text{KL}(p\|q) \geq 0$$

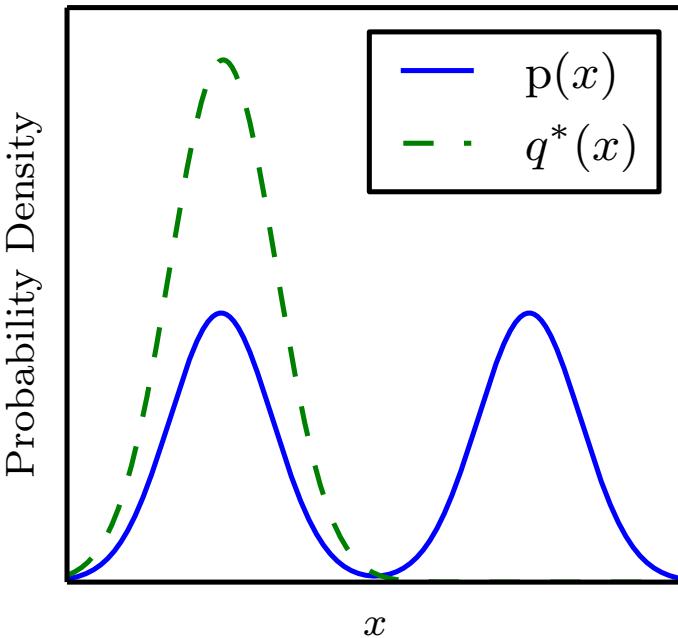
$$\text{KL}(p\|q) \neq \text{KL}(q\|p)$$

KL divergence is asymmetric!

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q\|p)$$



$$\text{KL}(p\|q) \not\equiv \text{KL}(q\|p)$$

Mutual Information

Quantifies how much information one variable carries about another

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned}$$

Symmetry: $I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$

Independent random variables: $p(x, y) = p(x)p(y) \Rightarrow I(x, y) = 0$

Questions?
