

Mathematics Preliminaries for Machine Learning CS 445/545



- Linear Algebra Overview
- Statistics/Probability Overview

# Linear Algebra: Vectors & Matrices

$$\mathbf{x} \in \mathbb{R}^d \to \mathbf{x} = \langle x_1, \dots, x_d \rangle$$

vector

$$\mathbf{A} \in \mathbb{R}^{m \times n} \rightarrow \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \xrightarrow{\mathbf{1d} \cdot \mathbf{ensor} \quad \mathbf{2d} \cdot \mathbf{ensor} \quad \mathbf{3d} \cdot \mathbf{ensor}}$$

$$\mathbf{M} \in \mathbb{R}^{m \times n \times p} \rightarrow \begin{bmatrix} a_{11}^{1} & \dots & a_{1n}^{1} \\ \vdots & \ddots & \vdots \\ a_{m1}^{1} & \dots & a_{mn}^{1} \end{bmatrix} \cdots \begin{bmatrix} a_{11}^{p} & \dots & a_{1n}^{p} \\ \vdots & \ddots & \vdots \\ a_{m1}^{p} & \dots & a_{mn}^{p} \end{bmatrix}$$

# Vector & Matrix Operations

dot product: 
$$\mathbf{x} \cdot \mathbf{y} = \sum_{i=1}^{d} x_i y_i \in \mathbb{R}$$

results in a scalar;  $\mathbf{x}$ ,  $\mathbf{y}$  are **orthogonal** if  $\mathbf{x} \cdot \mathbf{y} = 0$ 

matrix multiplication :  $A \cdot B = C$ 

 $m \times n \quad n \times p \quad m \times p$ 

$$A \cdot B = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \text{row i} & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & \dots & a_{1p} \\ \vdots & \vdots \\ b_{n1} & \dots & a_{np} \end{bmatrix} \rightarrow C_{ij} = \mathbf{a}_{rowi} \cdot \mathbf{b}_{colj}$$

• Matrix multiplication involves a sequence of dot products; element  $C_{ij}$  in the resultant matrix is equal to the dot product of *row i* (from the left matrix) and *column j* (from the right matrix).

# Vector & Matrix Operations

scalar multiplication:  $c\mathbf{x} = c \langle x_1, ..., x_d \rangle = \langle cx_1, ..., cx_d \rangle$ 

scalar multiplication: 
$$cA=c\begin{bmatrix}a_{11}&\ldots&a_{1n}\\\vdots&\ddots&\vdots\\a_{m1}&\ldots&a_{mn}\end{bmatrix}=\begin{bmatrix}ca_{11}&\ldots&ca_{1n}\\\vdots&\ddots&\vdots\\ca_{m1}&\ldots&ca_{mn}\end{bmatrix}$$

F 7

• Matrix multiplication is <u>associative</u>: A(BC) = (AB)C (always holds), but not <u>commutative</u>:  $AB \neq BA$  (in general).

vector transpose: 
$$\mathbf{x}_{1\times d} = [x_1, \dots, x_d] \rightarrow \mathbf{x}^T = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

matrix transpose: A=
$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \rightarrow A^{T} = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix}$$

• A matrix is called "symmetric" if:  $A^{T}=A$ .

# Vector Norms

• Norms convey the notion of the "magnitude" (i.e. size) of a vector; note that the equivalence (in magnitude) of two vectors is relative to the choice of norm.

• There are many types – even "families" – of norms relevant to ML/data science. Here are several of the most commonly used norms in ML:

(1) L2 norm (i.e. "Euclidean norm")

$$\|\mathbf{x}\|_{2} = \sqrt{\sum_{i=1}^{d} x_{i}^{2}}$$
 Note:  $\|\mathbf{x}\|_{2}^{2} = \sum_{i=1}^{d} x_{i}^{2} = \mathbf{x} \cdot \mathbf{x}$ 

(2) L1 norm (i.e. Manhattan distance)

$$\left\|\mathbf{x}\right\|_1 = \sum_{i=1}^d |x_i|$$

(3)  $\infty$  norm

$$\|\mathbf{x}\|_{\infty} = \max|x_i|$$

\*For an ML practitioner, the "choice" of a norm is oftentimes a crucial part of feature engineering and the ML problem formulation process itself; one can think of the different norm choices as striking a balance between "precision" and computational overhead.

\*There exist equivalent norms applied to matrices; the above norms are examples from the family of *p-norms*.



# Dot Product

• The dot product has important geometric properties that are useful in ML.

(\*) The dot product can be defined equivalently:

$$\mathbf{x} \cdot \mathbf{y} = \left\| \mathbf{x} \right\|_2 \left\| \mathbf{y} \right\|_2 \cos(\theta)$$



From this equivalent definition of the dot product, we can show that <u>the dot product</u> <u>quantifies the "similarity" between two vectors</u>. Consider (3) cases:

(i) Vectors **x** and **y** are "out of alignment" and meet at a 90 degree angle; in this case:

$$\mathbf{x} \cdot \mathbf{y} = \left\| \mathbf{x} \right\|_2 \left\| \mathbf{y} \right\|_2 \cos(90^\circ) = 0$$

(ii) Vectors **x** and **y** are "perfectly aligned" (i.e. *parallel* to one another):

 $\mathbf{x} \cdot \mathbf{y} = \left\| \mathbf{x} \right\|_{2} \left\| \mathbf{y} \right\|_{2} \cos \left( 0^{\circ} \right) = \left\| \mathbf{x} \right\|_{2} \left\| \mathbf{y} \right\|_{2}$ 

(iii) Vectors **x** and **y** are "oppositely aligned (i.e. they are *anti-parallel*):

 $\mathbf{x} \cdot \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos(180^\circ) = -\|\mathbf{x}\|_2 \|\mathbf{y}\|_2$ 

# Special Matrices

• The identity matrix I is a square ( $n \ge n$  matrix); the identity matrix multiplied by any matrix A (appropriately shaped) results in the matrix A:

$$I_{n} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & 1 \end{bmatrix}_{n \times n} \qquad AI = IA = A$$

• A matrix A is said to be **symmetric** if it equals its transpose:

$$A^{T} = A \qquad e.g., \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}^{T} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$$

Note: $(A+B)^{T} = A^{T} + B^{T}, (AB)^{T} = B^{T}A^{T}$ 

• For a **diagonal matrix**, all off-diagonal entries are zero (note that diagonal entries are permitted to be zero).

$$D = \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & d_m \end{bmatrix}$$

# Special Matrices

• An **upper-triangular matrix** has zero elements *below* the main diagonal. Note that Gaussian Elimination (from elementary linear algebra) yields an upper-triangular matrix.

| U = | <i>u</i> <sub>11</sub> | <i>u</i> <sub>12</sub> | ••• | $u_{1n}$          |
|-----|------------------------|------------------------|-----|-------------------|
|     | 0                      | <i>u</i> <sub>22</sub> | ••• | $u_{2n}$          |
|     | 0                      | 0                      | ·   | ÷                 |
|     | 0                      | 0                      |     | U <sub>mn</sub> _ |

• A lower-triangular matrix has zero elements above the main diagonal.

$$L = \begin{bmatrix} l_{11} & 0 & \cdots & 0 \\ l_{21} & l_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ l_{m1} & l_{m2} & \cdots & l_{mn} \end{bmatrix}$$

• An **orthogonal matrix** is a matrix with *orthonormal* rows and columns; equivalently, the <u>inverse of an orthogonal matrix is its transpose</u>.

$$Q^T Q = Q Q^T = I$$

# Special Matrices

• A square matrix *A* is **Positive Definite** if:

 $\forall \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}, \mathbf{x}^T A \mathbf{x} > 0$ 

• Analogously, a square matrix *A* is **Positive semi-Definite** if (e.g. covariance matrix):

## $\forall \mathbf{x} \in \mathbb{R}^n, \mathbf{x}^T A \mathbf{x} \ge 0$

• We say that the matrix  $A_{mxn}$  is **invertible** (i.e. *non-singular*) if there exists  $A^{-1}_{nxm}$ , where:

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

Properties:

$$(A^{T})^{-1} = A^{-T}$$
  $(AB)^{-1} = B^{-1}A^{-1}$  (if A,B non-singular)

# Linear Systems

• Commonly we encode, and subsequently solve systems of linear equations :

$$a_{m1}x_1 + \cdots + a_{mn}x_n = b_m$$

• When the *coefficient matrix* A is non-singular, the linear system gives rise to a unique solution:

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

\* Note that matrix inversion requires roughly on the order of  $O(n^3)$  arithmetic operations.

# Matrix Factorizations

• Matrix factorizations are immensely useful for identifying an underlying, inherent structure in a matrix (i.e. data).

Here are several important examples:

### LU Factorization

$$A = LU$$

• This factorization encodes the result of the Gaussian Elimination (GE) procedure (note that not all matrices admit of an LU factorization). L: denotes a lower-triangular matrix of "multipliers" used in GE. U denotes an upper-triangular (i.e. echelon form) matrix resulting from GE.

### PALU Factorization (Permuted LU factorization)

$$PA = LU$$

• This technique is similar to LU Factorization, except that we perform a **pivoting operation first** (i.e. permute the rows of *A* via a *permutation matrix*, *P*). LU factorization is subsequently performed; all matrices admit of such a factorization.

# Matrix Factorizations

Here are several important examples:

### **QR** Factorization

$$A = QR$$

• Q is an orthogonal matrix and R is upper-triangular -- commonly used for solving both regression problems and linear dynamical systems.

### **Eigendecomposition**

$$A = V \Sigma V^{T}$$

• This is one of the most useful and commonly-used of all matrix factorizations. The primary use of an eigendecomposition in ML is to perform *dimensionality reduction*; as such, this technique is closely related to **PCA** (*principal component analysis*) and **SVD** (*singular value decomposition* – see below) methods;  $\Sigma$  is a diagonal matrix consisting of the eigenvalues of A, and V is the matrix of corresponding eigenvectors.

# Matrix Factorizations

Here are several important examples:

**Cholesky Factorization** 

$$A = LL^T$$

• L is lower-triangular; Cholesky can be used to numerically solve linear systems; every positive-definite matrix admits of a Cholesky factorization.

### SVD (Singular Value Decomposition)

$$A = U\Sigma V^{T}$$

• SVD is one of the most essential matrix factorizations for applications of ML. U and V are orthogonal matrices, and  $\Sigma$  is a diagonal matrix containing the "singular values" (i.e. the eigenvalues of A<sup>T</sup>A. SVD has many applications (an orthogonal matrix denotes the matrix of eigenvectors of A<sup>T</sup>A, including dimensionality reduction and compression. All matrices admit of a singular value decomposition.

# Determinants

• Geometrically, the determinant of a square matrix A (written |A|) quantifies the unit increase in volume of the *linear transformation* defined by A (note that matrix multiplication defines a linear transformation).

Determinants can be computed through recursion; the general formula for a determinant is:

$$|A| = \sum_{i=1}^{n} a_{ij} (-1)^{i+j} \cdot M_{ij}$$
  
the "ij-minor"  
of A

Note that |A| = 0 if and only if A is singular.

Some Properties of Determinants:

$$|AB| = |A||B| \qquad |A^T| = |A|$$

# Eigenvalues

• The eigenvalues  $\lambda$  and eigenvectors **v** of a matrix A satisfy:

$$A\mathbf{v} = \lambda \mathbf{v} \qquad \left( v \neq 0 \right)$$

• Which means that the eigenvectors of a matrix A are precisely the vectors for which multiplication by A is tantamount to scalar multiplication by  $\lambda$ .



• Determining the exact values of the set of eigenvalues for a matrix  $A_{nxn}$  is requires solving the so-called **characteristic equation**:  $|A - \lambda I| = 0$ , which is an *n*-degree polynomial equation in the variable  $\lambda$ .

# Linear Independence, Span and Basis

• A set of vectors is called linearly independent if the set contains "no redundancy"; formally:

**Def**. A set of vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is **linearly independent** if:  $\alpha_1 \mathbf{v}_1 + \dots + \alpha_n \mathbf{v}_n = \mathbf{0}$  implies  $\alpha_1 = 0 \forall 1 \le i \le n$ 

• The span of a set is defined as the set of all linear combinations of the set of vectors.

**Def**. span 
$$\{\mathbf{v}_1, \dots, \mathbf{v}_n\} = \{\alpha_{i1}\mathbf{v}_1 + \cdots + \alpha_{in}\mathbf{v}_n \mid i \in \mathbb{R}\}$$

• A basis is a set of linearly independent vectors that spans the parent vector space.

$$e.g., \left\{ \begin{bmatrix} 1\\0\\0 \end{bmatrix}, \begin{bmatrix} 0\\1\\0 \end{bmatrix}, \begin{bmatrix} 0\\0\\1 \end{bmatrix} \right\} \text{ is the "standard" basis set for } \mathbb{R}^3$$

## The Four Fundamental Subspaces

- The Four Fundamental subspaces of a matrix A<sub>mxn</sub>:
  - 1. (Column Space) Col(A): the span of the column vectors of A.
  - 2. (Null Space) Null(A): the set of all vectors that satisfy Ax=0.
  - 3. (Row Space) Row(A): the span of the row vectors of A.
  - 4. (Null Space of  $A^{T}$ ) Null( $A^{T}$ ): the set of all vectors that satisfy  $A^{T}x=0$ .



• We use statistics and probability to quantify and summarize our beliefs about a "state of the world" in the face of incomplete or partial knowledge.

• Denote a random event E; the **sample space S** consists of the set of all possible outcomes associated with E (e.g. if E="coin flip",  $S=\{H,T\}$ ).

• A **random variable** (e.g. X, Y) is a variable that is assigned a number based on the outcome of the random event E.

• Random variables are either **Discrete** (e.g., 0/1) or **Continuous** (e.g., height, time).

### Overview of Statistics/Probability <u>Probability Distributions</u>

• A **probability distribution** summarizes our total knowledge about the random event E, via the random variable X.

• For a discrete random variable, the probability distribution of X is called a **probability mass function** (pmf); a pmf satisfies the following properties, with |S|=k:

1.  $0 \le p(X_i) \le 1$   $1 \le i \le k$ 2.  $\sum_{i=1}^k p(X_i) = 1$ 

• Similarly, for a continuous random variable, the probability distribution of X is called a **probability density function** (pdf); a pdf satisfies the following properties, with  $|S| = \infty$ :

$$1. \ 0 \le p(X_i) \le 1$$
$$2. \int_{S} p(X) dx = 1$$

### Overview of Statistics/Probability <u>Probability Distributions</u>

• A **cumulative density function** (cdf) is defined as the cumulative probability up to a given value of a random variable:

$$F_{X}(x) = p(X \le x) = \int_{-\infty}^{x} p(u) \, du$$

• Percentiles and quartiles can be defined in a natural way with respect to a cdf:

$$\underbrace{F_X(x) = 0.25}_{x=Q_1} \qquad \underbrace{F_X(x) = 0.5}_{x=Q_2 \text{ (median)}} \qquad \underbrace{F_X(x) = 0.75}_{x=Q_3}$$

• Note that due to the Fundamental Theorem of Calculus, it follows that:

$$\frac{d}{dx} \underbrace{F_{X}(x)}_{cdf} = p(x)$$



### **Properties of Probability Distributions**

- Two random events  $E_1$  and  $E_2$  are disjoint if:  $S_1 \cap S_2 = \emptyset$ .
- If two events  $E_1$  and  $E_2$  are <u>disjoint</u>, then:  $P(E_1 \cup E_2) = P(E_1) + P(E_2)$

• More generally, the addition rule of probability states, that for *any* two events  $E_1$  and  $E_2$ :

$$P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1 \cap E_2)$$

### **Conditional Probability**

• Def. Conditional Probability:

 $P(A | B) = \frac{P(A \cap B)}{P(B)}$ "A given B'

• From this definition, we can derive the multiplication rule of probability:

$$P(A \cap B) = P(A \mid B)P(B)$$

• Equivalently,

$$P(A \cap B) = P(B \mid A)P(A)$$

• We say that events A & B are **independent** if the outcome of A has no bearing on B (and vice versa); more formally the joint probability distribution p(A,B) factors.

• Def. A & B are independent if:

$$P(A \cap B) = P(A)P(B)$$

• Equivalently, if A & B are independent, it also follows that:

$$P(A | B) = P(A)$$
  $P(B | A) = P(B)$ 

Thus, in summary, if A & B are independent:  $P(A \cap B) = (A | B)P(B) = P(A)P(B)$ 

\*Independence is commonly denoted:  $A \perp B$ 

• Two major theorems in elementary statistics: (1) the Law of Large Numbers and (2) the Central Limit Theorem.

• The Law of Large Numbers states (paraphrasing): Experimental (i.e. empirical) probabilities converge to their associate theoretical probability as the number of trials tends to infinity.



The Central Limit Theorem (a conceptual pillar of statistics)

**In words**: given a sufficiently large sample size from a population (with a finite level of variance), the mean of all samples from the same population will be approximately equal to the mean of the population. Furthermore, all of the samples will follow an approximate normal distribution pattern, with all variances being approximately equal to the variance of the population divided by each sample's size.



In a picture:

Whatever the form of the population distribution, the sampling distribution tends to a Gaussian, and its dispersion is given by the Central Limit Theorem.

<u>In a theorem</u>: Suppose  $\{X_1, X_2, ...,\}$  is a sequence of I.I.D. random variables with  $E[X_i] = \mu$  and  $Var[X_i] = \sigma^2 < \infty$ Then as n approaches infinity, the random variable  $(1/n)(X_1 + ... + X_n)$  converges to a normal  $N(0, \sigma^2/n)$ :

$$\left(\left(\frac{1}{n}\sum_{i=1}^{n}X_{i}\right)-\mu\right)^{d} \rightarrow N\left(0,\frac{\sigma^{2}}{n}\right)$$

## Probability Distributions

• Here are some (but certainly not all) of the essential probability distributions for ML and applied statistics:

### 1-D Gaussian (i.e. Normal)

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



• When  $\mu=0$  and  $\sigma=1$  (i.e. N(0,1)) we call this the standard Normal model.

## Probability Distributions

### Multivariate Gaussian (i.e. MVN)

$$N(x;\mu,\Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$



Using different forms of covariance matrix allows for clusters of different shapes



## Probability Distributions

### Bernoulli & Binomial Distributions

• The **Bernoulli distribution** is a single variable, discrete distribution, describing a random variable with two discrete states (e.g. heads/tails for a single coin flip). The Bernoulli distribution forms the basis of the **Binomial distribution**, which models repetitions of independent Bernoulli trials (N total).

### Bernoulli pmf

$$p(X=1) = \theta$$
  $p(X=0) = 1 - \theta$   $0 \le \theta \le 1$ 

### **Binomial pmf**

$$\underbrace{p(X=k)}_{\text{probability of successes in n trials"}} = \binom{n}{k} \theta^k (1-\theta)^{n-k}$$

 $p(\text{exactly 7 H in 10 flips}) = {\binom{10}{7}} 0.6^7 (1-0.6)^{10-7}$ 

• For example, with a biased coin ( $\theta = 0.6$ ), we have:

# Summary Statistics for Random Variables

### Expectation and Variance of a Random Variable

• The **Expected Value** of a random variable X summarizes the outcome: "if the trial were executed once, on average, this is the numerical value we would expect for X"; E[X] accordingly computes the <u>arithmetic mean of a random variable</u>, i.e.  $E[X]=\mu$ .

$$E[X] = \sum_{i=1}^{k} x_i P(X = x_i) \quad \text{(Discrete RV)} \quad E[X] = \int_{S} x p(x) dx \quad \text{(Continuous RV)}$$

For example, to compute the expected number of heads X in 10 flips of a fair coin (X~Binomial) we have:

$$E[X] = \sum_{i=0}^{10} k \binom{10}{k} \left(\frac{1}{2}\right)^k \left(\frac{1}{2}\right)^{n-k} = 5$$

### Summary Statistics for Random Variables

### **Properties of Expected Value and Variance**

• Expected Value is a **linear operator** (as are matrix multiplication, limits, differentiation and integration, among other common mathematical operators) -- meaning that it obeys the following two *linearity properties:* 

1.For any two random variables X, Y: E[X+Y] = E[X] + E[Y]2.For any  $c \in \mathbb{R}$  : E[cX] = cE[X]

The following corollary is also useful: 
$$Var[X] = E[X^2] - E[X]^2$$

### Proof.

$$Var[X] = E[(X - \mu)^{2}] = E[X^{2} - 2\mu X + \mu^{2}]$$
  

$$= E[X^{2}] - E[2\mu X] + E[\mu^{2}]$$
  
by linearity  
of E[·]  

$$= E[X^{2}] - 2\mu E[X] + \mu^{2}$$
  
by linearity  
of E[·]  

$$= E[X^{2}] - 2\mu^{2} + \mu^{2}$$

### Covariance

• **Covariance** is a measure of the <u>linear relationship</u> between two random variables, X and Y. If Cov(X,Y) > 0, this indicates a positive linear relationship between the random variables (i.e. as X increases, Y increases; as X decreases, Y decreases); when Cov(X,Y) < 0 the variables share a negative linear relationship; Cov(X,Y) = 0 indicates the absence of a linear relationship.

Def. 
$$Cov(X,Y) = E\left[(X-E[Y])(X-E[Y])\right]$$

#### Lemma

If  $X \perp Y$  (*i.e.* if X and Y are independent), then Cov(X,Y) = 0

\* Note that the *converse* of the lemma above fails; in other words Cov(X,Y)=0 need not imply that X and Y are independent.

### Covariance

• The Covariance Matrix ( $\Sigma$ ) for a set of random variables {X<sub>1</sub>,...,X<sub>N</sub>} is defined as the matrix of pairwise covariances:

Def.

Let 
$$\mathbf{X} = \{X_1, ..., X_N\}, \quad \Sigma_{ij} = Cov(X_i, X_j) = E\left[\left(X_i - E\left[X_i\right]\right)\left(X_j - E\left[X_j\right]\right)\right]$$
  
matrix of  
covariances

$$\Sigma = \begin{bmatrix} Var[X_1] & Cov(X_1, X_2) & \cdots & Cov(X_1, X_N) \\ Cov(X_2, X_1) & Var[X_2] & \cdots & Cov(X_2, X_N) \\ \vdots & \vdots & \ddots & \vdots \\ Cov(X_N, X_1) & Cov(X_N, X_2) & \cdots & Var[X_N] \end{bmatrix}$$

Note that  $\Sigma$  is symmetric and positive semi-definite. The covariance matrix is used to parameterize the MVN (multivariate normal distribution); the covariance matrix can likewise be computed for a dataset.

## Bayes' Theorem

• **Bayes' Theorem** is a vital (yet simple) conditional probability formula; today its use is omnipresent across ML.

Def. 
$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$
  
Derivation: 
$$P(A | B) = \underbrace{\frac{P(A \cap B)}{P(B)}}_{\text{by definition of conditional probability}} = \underbrace{\frac{P(B | A)P(A)}{P(B)}}_{\text{by multiplication rule}}$$

• More importantly, <u>Bayes' Theorem can be generalized to encapsulate the whole of the inductive element of the scientific method</u>. To this end, consider H (hypothesis) and D (data):

$$P(H \mid D) = \frac{P(D \mid H)P(H)}{P(D)}$$

• In this case, Bayes' Theorem yields a natural mechanism for <u>updating our belief about the</u> <u>world/the plausibility of a hypothesis (H) given an observation (D)</u>. P(H|D) is referred to as the **posterior probability** of H, P(H) is called the **prior probability** of H, P(D|H) defines the **likelihood of the data**, and P(D) is the **data prior**.

## Bayesian and Frequentist Statistics

• There exist two general paradigms for modern statistics: the **frequentist** and **Bayesian** approaches.

**Frequentists:** Generally consider model parameters ( $\theta$ ) as fixed; data are drawn from some objective distribution, defined by  $\theta$ . There exists various well-known pathologies associated with frequentism, including the "problem of induction" (Hume), the Black Swan Paradox, limited exact solutions and a heavy reliance upon long-term frequencies.

**Bayesians**: (Observed) data are fixed; data are observed from a realized sample; we encode prior beliefs, and parameters values are described probabilistically.

• Frequentists use the Maximum Likelihood Estimate (MLE) for point estimates of parameters  $\theta$ :

$$\hat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$

• Bayesians instead use the Maximum A Posterior (MAP) for parameter estimates:

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta \mid D) = \arg \max_{\theta} P(D \mid \theta) P(\theta)$$

• The **entropy** of a discrete random variable X (equivalently: the entropy of the pmf associated with X) is defined:

$$H(X) = -\sum_{i} p(X = x_i) \log p(X = x_i)$$

• The differential entropy of a continuous random variable is defined analogously:

$$H(X) = -\int_{S} p(x) \log p(x) dx$$

• Entropy quantifies disorder/"surprise"; the **Principle of Insufficient Reasons** (PIR) states (paraphrasing) that in the absence of compelling evidence, one should adopt a **maximum entropy probability distribution**. The uniform distribution is a maximum entropy distribution; the Gaussian distribution is a likewise a maximum entropy distribution (up to second moments). Entropy is minimized (i.e. zero) for deterministic events, e.g. *Dirac delta function*.



Ex. The entropy of a Bernoulli random variable X is given by:

$$p(X=1) = \theta \quad p(X=0) = 1 - \theta \quad 0 \le \theta \le 1$$
$$H(X) = -\sum_{i} p(X=x_{i}) \log p(X=x_{i})$$
$$= -(p(X=1) \log p(X=1) + p(X=0) \log p(X=0))$$
$$= -(\theta \log \theta + (1-\theta) \log (1-\theta))$$



\* Notice that entropy is maximized in this case when  $\theta = 0.5$ , which corresponds with a binary uniform distribution; conversely, entropy is minimized when either  $\theta = 0$  or  $\theta = 1$ , in which case the even is deterministic.

• The Kullback-Leibler Divergence quantifies the <u>difference between two probability</u> <u>distributions</u>, p(x) and q(x).

Def.

$$KL(p \parallel q) = \sum_{i} p(X = x_{i}) \log \frac{p(X = x_{i})}{q(X = x_{i})}$$
$$KL(p \parallel q) = \int_{S} p(x) \log \frac{p(x)}{q(x)} dx$$



The Information Inequality states:

$$KL(p \parallel q) \ge 0$$
 and  $KL(p \parallel q) = 0 \leftrightarrow p = q$ 

$$KL(p || q) = \sum_{i} p(X = x_{i}) \log \frac{p(X = x_{i})}{q(X = x_{i})}$$
$$KL(p || q) = \int_{S} p(x) \log \frac{p(x)}{q(x)} dx$$

$$KL(p || q) \ge 0$$
 and  $KL(p || q) = 0 \leftrightarrow p = q$ 

\*Recall that covariance/correlation are inherent measures of the <u>linear</u> relationship between two random variables. Using KL-divergence, we can develop a more general notion of independence, called **mutual information**.

Def.

$$MI(X,Y) = KL(p(X,Y) || p(X)p(Y)) = \sum_{X} \sum_{Y} p(X,Y) \log \frac{p(X,Y)}{p(X)p(Y)}$$
$$MI(X,Y) = KL(p(X,Y) || p(X)p(Y)) = \int_{S_X} \int_{S_Y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} dxdy$$

\*From the information inequality, it follows that:

$$MI(X,Y) \ge 0$$
 and  $MI(X,Y) = 0 \leftrightarrow p(X,Y) \parallel p(X)p(Y)$   
(i.e.  $MI(X,Y) = 0 \leftrightarrow X \perp Y$ )

Thus, MI can be seen as a more general measure of statistical independence than covariance.



