

CS 559: Machine Learning Fundamentals & Applications

Lecture 2: Mathematics for Machine Learning

Spring 2023



Outline

- Linear Algebra
- Analytic Geometry
- Vector Calculus
- Probability Theory





2.1. Linear Algebra



2.1.1. Basic Matrix Identities

A matrix \mathbf{A} has elements A_{ij} where i indexes the rows, and j indexes the columns.

$$\mathbf{A} = \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix}, \quad A_{ij} \in \mathbb{R}.$$



2.1.2. Matrix Addition and Multiplication

The sum of two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{m \times n}$ is the element-wise sum,

$$\mathbf{C} = \mathbf{A} + \mathbf{B} = \begin{bmatrix} a_{11} + b_{11} & \cdots & a_{1n} + b_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} + b_{m1} & \cdots & a_{mn} + b_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}.$$

The product of two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times k}$ is the pair-wise sum,

$$\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times k}, c_{ij} = \sum_{l=1}^n a_{il} b_{lj}, i = 1, \dots, m, j = 1, \dots, k.$$



2.1.2. Matrix Addition and Multiplication

Matrix Multiplication Properties:

- Not commutative: $\mathbf{AB} \neq \mathbf{BA}$
- Associative: $\forall \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times p}, \mathbf{C} \in \mathbb{R}^{p \times q}: (\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$
 - $(\lambda\psi)\mathbf{C} = \lambda(\psi)\mathbf{C}$ where λ and ψ are constants.
 - $\psi(\mathbf{BC}) = (\psi\mathbf{B})\mathbf{C} = \mathbf{B}(\psi\mathbf{C}) = (\mathbf{BC})\psi, \mathbf{B} \in \mathbb{R}^{m \times n}, \mathbf{C} \in \mathbb{R}^{n \times k}$
- Distributive: $\forall \mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}, \mathbf{C} \in \mathbb{R}^{n \times p}: (\mathbf{A} + \mathbf{B})\mathbf{C} = \mathbf{AC} + \mathbf{BC}$
 - $(\lambda + \psi)\mathbf{C} = \lambda\mathbf{C} + \psi\mathbf{C}$



2.1.3. Inverse and Transpose

A matrix \mathbf{I}_N is the $N \times N$ *identity* matrix (also called the unit matrix)

$$\mathbf{I}_N = \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{bmatrix}.$$

The inverse of $\mathbf{A} \in \mathbb{R}^{n \times n}$, \mathbf{A}^{-1} , satisfies

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}.$$

If the inverse exists, \mathbf{A} , is called *regular/invertible/nonsingular*.

Consider two matrices

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \text{ and } \mathbf{A}' = \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix},$$

the product of $\mathbf{A}\mathbf{A}'$ is then

$$\mathbf{A}\mathbf{A}' = \begin{bmatrix} a_{11}a_{22} - a_{12}a_{21} & 0 \\ 0 & a_{11}a_{22} - a_{12}a_{21} \end{bmatrix} = (a_{11}a_{22} - a_{12}a_{21})\mathbf{I}.$$

If and only if $a_{11}a_{22} - a_{12}a_{21} \neq 0$,

$$\mathbf{A}^{-1} = \frac{1}{a_{11}a_{22} - a_{12}a_{21}} \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}.$$



2.1.3. Inverse and Transpose

Inverse Matrix Properties:

- $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I} = \mathbf{A}^{-1}\mathbf{A}$
- $(\mathbf{A}\mathbf{B})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$
- $(\mathbf{A} + \mathbf{B})^{-1} \neq \mathbf{A}^{-1} + \mathbf{B}^{-1}$

The *Woodbury identity* is

$$(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}.$$

The transpose matrix \mathbf{A}^T has elements $(\mathbf{A}^T)_{ij} = A_{ji}$. From the definition of transpose, a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *symmetric* if $\mathbf{A} = \mathbf{A}^T$.

Transpose Properties:

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

If \mathbf{A} is invertible, then $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1} = \mathbf{A}^{-T}$.



2.1.4. Vector Spaces

Vector: Each row or column in a matrix \mathbf{A} is called a *vector*. A real-valued vector $V = (\mathcal{V}, +, \cdot)$ is a set \mathcal{V} with two operations

$$\begin{aligned} +: \mathcal{V} \times \mathcal{V} &\rightarrow \mathcal{V} \\ \cdot: \mathbb{R} \times \mathcal{V} &\rightarrow \mathcal{V} \end{aligned}$$

where $+$ is the vector addition and \cdot is a multiplication by a scalar.

Vector Subspace: Suppose V is a vector space and $\mathcal{U} \subseteq \mathcal{V}, \mathcal{U} \neq \emptyset$. Then U is a vector subspace of V if U is a **vector subspace** with the vector space operations restricted to $\mathcal{U} \times \mathcal{U}$ and $\mathbb{R} \times \mathcal{U}$.

Linear Combination: A vector space V and a finite number of vectors $\mathbf{x}_1, \dots, \mathbf{x}_k \in V$. Then, every $\mathbf{v} \in V$ of the form

$$\mathbf{v} = \sum_{i=1}^k \lambda_i \mathbf{x}_i \in V$$

with $\lambda_i \in \mathbb{R}$ is a *linear combination* of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$.



2.1.5. Basis and Rank

Span: When a set of vectors $\mathcal{A} = \{x_1, \dots, x_k\}$ is in a vector space V and if every vector can be expressible in a linear combination format, \mathcal{A} is called the *generating set* of V . The set of all linear combinations of vectors in \mathcal{A} is called the **span** of \mathcal{A} . The denotation $V = \text{span}[\mathcal{A}]$ means \mathcal{A} spans the vector space V .

Basis: If there exists no smaller set $\mathcal{A} \subseteq V$ that spans, every linearly independent generating set of V is called a **basis** of V .

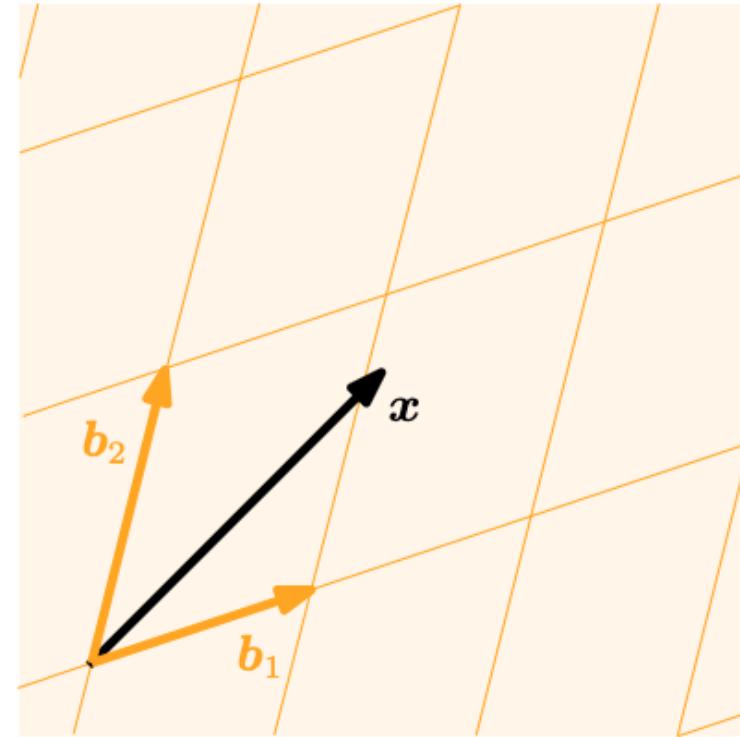
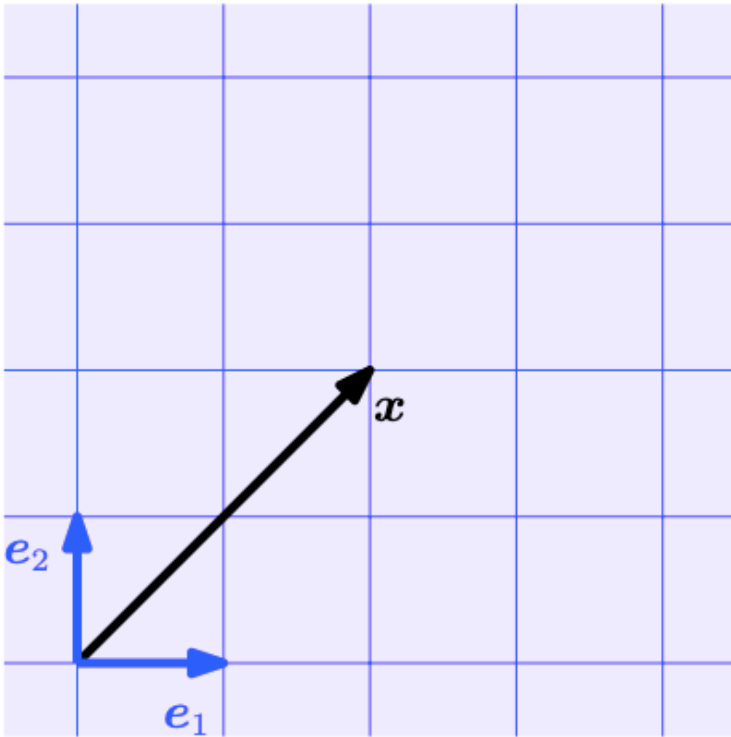
Rank: The number of linear independent columns of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ equals the number of linearly independent rows and is called the **rank** of \mathbf{A} .

2.1.6. Linear Mapping

Linear Mapping: For vector spaces V, W , a mapping $\Phi: V \rightarrow W$ is called a **linear mapping** if

$$\forall \mathbf{x}, \mathbf{y} \in V, \forall \lambda, \psi \in \mathbb{R}: \Phi(\lambda \mathbf{x} + \psi \mathbf{y}) = \lambda \Phi(\mathbf{x}) + \psi \Phi(\mathbf{y}).$$

- If $\Phi: V \rightarrow W, \Psi: V \rightarrow W$ are linear, then $\Phi + \Psi$ and $\lambda \Phi, \lambda \in \mathbb{R}$, are linear, too.





2.1.6. Linear Mapping

In a vector space V , an ordered basis $B = (\mathbf{b}_1, \dots, \mathbf{b}_n)$ of V , a unique linear combination of $\mathbf{x} \in V$ is formulated:

$$\mathbf{x} = \alpha_1 \mathbf{b}_1 + \dots + \alpha_n \mathbf{b}_n$$

where $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_n] \in \mathbb{R}^n$ are coordinates of \mathbf{x} w.r.t. B .

Consider vector spaces V, W with corresponding bases $B = (b_1, \dots, b_n)$ and $C = (c_1, \dots, c_m)$. The linear mapping $\Phi: V \rightarrow W$ for $j \in \{1, \dots, n\}$ is

$$\Phi(\mathbf{b}_j) = \sum_{i=1}^m \alpha_{ij} \mathbf{c}_i.$$

The $m \times n$ matrix \mathbf{A}_{Φ} is a transformation matrix of Φ .



2.2. Analytic Geometry



2.2.2. Norms

Norm: A **norm** on a vector space V is a function

$$\begin{aligned} ||\cdot||: V &\rightarrow \mathbb{R}, \\ \mathbf{x} &\mapsto ||\mathbf{x}||, \end{aligned}$$

which assigns each vector \mathbf{x} its *length* $||\mathbf{x}|| \in \mathbb{R}$, such that for all $\lambda \in \mathbb{R}$ and $\mathbf{x}, \mathbf{y} \in V$ the following hold:

- *Absolutely homogeneous:* $||\lambda\mathbf{x}|| = |\lambda| ||\mathbf{x}||$
- *Triangle inequality:* $||\mathbf{x} + \mathbf{y}|| \leq ||\mathbf{x}|| + ||\mathbf{y}||$
- *Positive definite:* $||\mathbf{x}|| \geq 0$ and $||\mathbf{x}|| = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}$.

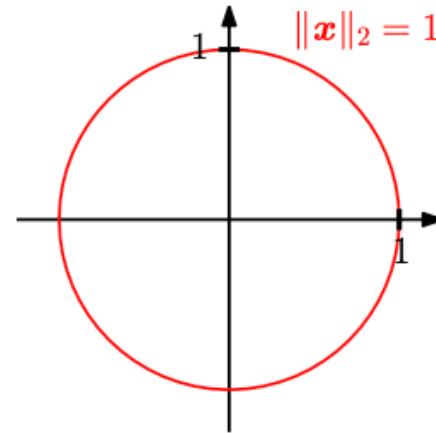
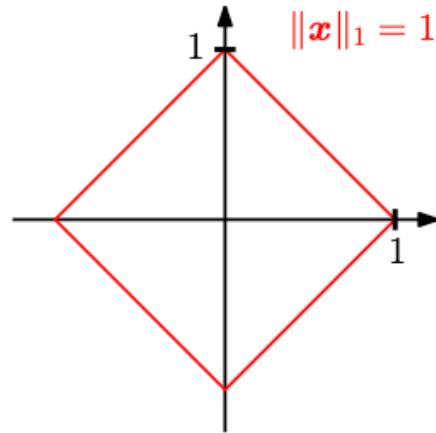
2.2.2. Norms

Manhattan Norm: The **Manhattan norm** (also called l_1) on \mathbb{R}^n is defined for $\mathbf{x} \in \mathbb{R}^n$ as

$$\|\mathbf{x}\| = \sum_{i=1}^n |x_i|.$$

Euclidean Norm: The **Euclidean norm** (also called l_2) of $\mathbf{x} \in \mathbb{R}^n$ is defined as

$$\|\mathbf{x}\| = \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{\mathbf{x}^T \mathbf{x}}.$$





2.2.2. Inner Product

Dot Product: The dot product takes two equal-length sequences of numbers or vectors and returns a single number as follow

$$\mathbf{x}^T \mathbf{y} = \mathbf{x}^T \cdot \mathbf{y} = \sum_{i=1}^n x_i y_i.$$

A *bilinear mapping* Ω is a mapping with two arguments, and it is linear in each argument that holds for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V, \lambda, \psi \in \mathbb{R}$ such that

$$\Omega(\lambda \mathbf{x} + \psi \mathbf{y}, \mathbf{z}) = \lambda \Omega(\mathbf{x}, \mathbf{z}) + \psi \Omega(\mathbf{y}, \mathbf{z})$$

$$\Omega(\mathbf{x}, \lambda \mathbf{y} + \psi \mathbf{z}) = \lambda \Omega(\mathbf{x}, \mathbf{y}) + \psi \Omega(\mathbf{x}, \mathbf{z})$$

- Ω is called *symmetric* if $\Omega(\mathbf{x}, \mathbf{y}) = \Omega(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in V$, i.e., the order of the argument does not matter.
- Ω is called *positive definite* if

$$\forall \mathbf{x} \in V \setminus \{\mathbf{0}\} : \Omega(\mathbf{x}, \mathbf{x}) > 0, \Omega(\mathbf{0}, \mathbf{0}) = 0.$$

- A positive definite, symmetric bilinear mapping $\Omega: V \times V \rightarrow \mathbb{R}$ is an *inner product* on V .
- The denotation of inner product is $\langle \mathbf{x}, \mathbf{y} \rangle$.
- An inner product is not always the dot product! For example, we define the inner product in \mathbb{R}^2 is

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 - (x_1 y_2 + x_2 y_1) + 2x_2 y_2,$$

it is different from $\mathbf{x}^T \mathbf{y} = x_1 y_1 + x_2 y_2$.

- The inner product of two vectors

$$\sqrt{\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle} = \|\mathbf{x} - \mathbf{y}\|$$

is the *distance* between two vectors.



2.2.2. Inner Product

Consider an n -dimensional vector space V with an inner product and an ordered basis $B = (\mathbf{b}_1, \dots, \mathbf{b}_n)$ of V . Any vectors $\mathbf{x}, \mathbf{y} \in V$ form as linear combinations $\mathbf{x} = \sum_{i=1}^n \psi_i \mathbf{b}_i \in V$ and $\mathbf{y} = \sum_{j=1}^n \lambda_j \mathbf{b}_j \in V$ where $\forall \lambda_j, \psi_i \in \mathbb{R}$. The inner product is

$$\forall \mathbf{x}, \mathbf{y}: \langle \mathbf{x}, \mathbf{y} \rangle = \left\langle \sum_{i=1}^n \psi_i \mathbf{b}_i, \sum_{j=1}^n \lambda_j \mathbf{b}_j \right\rangle = \sum_{i=1}^n \sum_{j=1}^n \psi_i \langle \mathbf{b}_i, \mathbf{b}_j \rangle \lambda_j = \hat{\mathbf{x}}^T \mathbf{A} \hat{\mathbf{y}},$$

where $A_{ij} = \langle \mathbf{b}_i, \mathbf{b}_j \rangle$ and $\hat{\mathbf{x}}, \hat{\mathbf{y}}$ are unit vectors $\left(\hat{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|}, \hat{\mathbf{y}} = \frac{\mathbf{y}}{\|\mathbf{y}\|} \right)$ with respect to B .

A symmetric matrix that satisfies

$$\forall \mathbf{x} \in V \setminus \{\mathbf{0}\}: \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$$

is called *symmetric, positive semidefinite*.



2.2.3. Orthogonality

Orthogonality: Two vectors \mathbf{x} and \mathbf{y} are **orthogonal** if and only if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ and its denotation is $\mathbf{x} \perp \mathbf{y}$. That is the angle between \mathbf{x} and \mathbf{y} is 0,

$$\cos \omega = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\sqrt{\langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle}} = \frac{\mathbf{x}^T \mathbf{y}}{\sqrt{\mathbf{x}^T \mathbf{x} \mathbf{y}^T \mathbf{y}}} = 0.$$

A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is an *orthogonal matrix* i.f.f. its columns are orthogonal so that

$$\mathbf{A} \mathbf{A}^T = \mathbf{I} = \mathbf{A}^T \mathbf{A},$$

which implies that

$$\mathbf{A}^{-1} = \mathbf{A}^T.$$

Orthonormal Basis: The basis is called an **orthonormal basis** if an n -dimensional vector and a basis $\{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ hold

$$\langle \mathbf{b}_i, \mathbf{b}_j \rangle = 0 \text{ for } i \neq j \text{ and } \langle \mathbf{b}_i, \mathbf{b}_i \rangle = 1$$

for all $i, j = 1, \dots, n$.

Projection: A linear mapping $\pi: V \rightarrow U$ is called a **projection** if $\pi^2 = \pi \circ \pi = \pi$ where $U \subseteq V$. The *projection matrices* \mathbf{P}_π exhibit the property that $\mathbf{P}_\pi^2 = \mathbf{P}_\pi$.

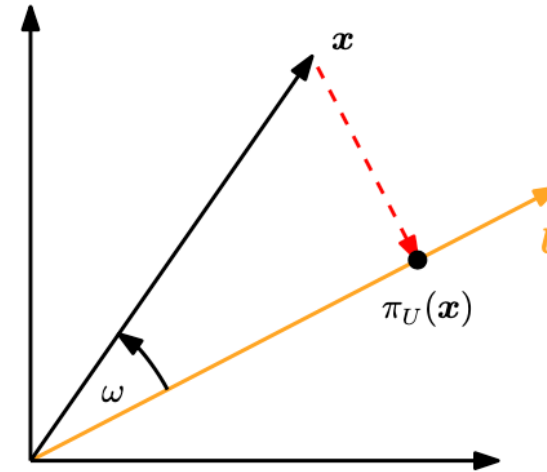
2.2.3. Orthogonality

- The segment $\pi_U(\mathbf{x}) - \mathbf{x}$ is orthogonal to U and therefore the basis vector \mathbf{b} of U , $\langle \pi_U(\mathbf{x}) - \mathbf{x}, \mathbf{b} \rangle = 0$.
- The projection $\pi_U(\mathbf{x})$ of \mathbf{x} onto U must be an element of U and, therefore, \mathbf{b} spans U , $\pi_U(\mathbf{x}) = \lambda \mathbf{b}$ where λ is the coordinate

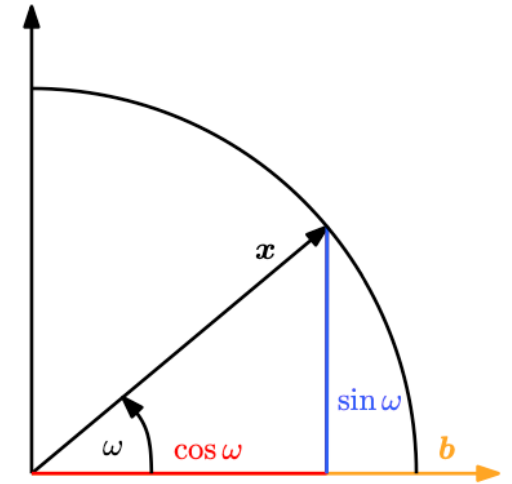
$$\lambda = \frac{\mathbf{b}^T \mathbf{x}}{\mathbf{b}^T \mathbf{b}} = \frac{\mathbf{b}^T \mathbf{x}}{\|\mathbf{b}\|^2}$$

and

$$\pi_U(\mathbf{x}) = \lambda \mathbf{b} = \frac{\mathbf{b}^T \mathbf{x}}{\|\mathbf{b}\|^2} \mathbf{b}.$$



(a) Projection of $\mathbf{x} \in \mathbb{R}^2$ onto a subspace U with basis vector \mathbf{b} .



(b) Projection of a two-dimensional vector \mathbf{x} with $\|\mathbf{x}\| = 1$ onto a one-dimensional subspace spanned by \mathbf{b} .



2.2.4. Determinant and Trace

Determinant: A **determinant** is a mathematical object in the analysis and solution of systems of linear equations. It is only defined for *square matrices* and it maps a square matrix onto a real number. Recall the inverse matrix of $\mathbf{A} \in \mathbb{R}^{2 \times 2}$, \mathbf{A}^{-1} ,

$$\mathbf{A}^{-1} = \frac{1}{a_{11}a_{22} - a_{12}a_{21}} \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}$$

can be expressed using the determination as following

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}$$

where

$$\det(\mathbf{A}) = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21}.$$

If $n = 3$,

$$\begin{aligned} \det(\mathbf{A}) &= \begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = a_{11}(a_{22}a_{33} - a_{23}a_{32}) \\ &\quad - a_{12}(a_{21}a_{33} - a_{23}a_{31}) + a_{13}(a_{21}a_{32} - a_{22}a_{31}) \end{aligned}$$



2.2.4. Determinant and Trace

For any $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\det(\mathbf{A})$ also can be computed as

- Expansion along column j :

$$\det(\mathbf{A}) = \sum_{k=1}^n (-1)^{k+j} a_{kj} \det(\mathbf{A}_{kj}).$$

- Expansion along row j :

$$\det(\mathbf{A}) = \sum_{k=1}^n (-1)^{k+j} a_{jk} \det(\mathbf{A}_{jk}).$$

where $\mathbf{A}_{kj} \in \mathbb{R}^{(n-1) \times (n-1)}$ is the *submatrix* of \mathbf{A} that is obtained when row k and column j are deleted.

For example,

$$\begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = (-1)^{1+1} a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} + (-1)^{1+2} a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + (-1)^{1+3} a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}.$$

The determinant exhibits the following properties:

- $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$
- If \mathbf{A} is invertible, then $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$
- Adding a multiple of a column/row to another one does not change $\det(\mathbf{A})$.
- Multiplication of a column/row with $\lambda \in \mathbb{R}$ scales $\det(\mathbf{A})$ by λ : $\det(\mathbf{A}) = \lambda^n \det(\mathbf{A})$.
- Swapping two rows/columns changes the sign of $\det(\mathbf{A})$.



2.2.4. Determinant and Trace

Trace: The **trace** of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is defined as

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii}$$

and satisfies the following properties:

- $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$ for $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$.
- $\text{tr}(\alpha \mathbf{A}) = \alpha \text{tr}(\mathbf{A})$, $\alpha \in \mathbb{R}$ for $\mathbf{A} \in \mathbb{R}^{n \times n}$.
- $\text{tr}(\mathbf{I}_n) = n$.
- $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ for $\mathbf{A} \in \mathbb{R}^{n \times k}$, $\mathbf{B} \in \mathbb{R}^{k \times n}$.
- For vectors, \mathbf{x}, \mathbf{y} , $\text{tr}(\mathbf{xy}^T) = \mathbf{y}^T \mathbf{x} \in \mathbb{R}$.



2.2.5. Eigenvalues and Eigenvectors

Eigenvalues and Eigenvector: Let $A \in \mathbb{R}^{n \times n}$ be a square matrix. Then $\lambda \in \mathbb{R}$ is an **eigenvalue** of A and $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ is the corresponding **eigenvector** of A if

$$A\mathbf{x} = \lambda\mathbf{x}$$

where λ is eigenvalue of $A \in \mathbb{R}^{n \times n}$. It is equivalent to

$$(A - \lambda I_n) = 0$$

and therefore

$$\det(A - \lambda I_n) = 0.$$



2.3. Vector Calculus

2.3.2. Differentiation of Univariate Functions

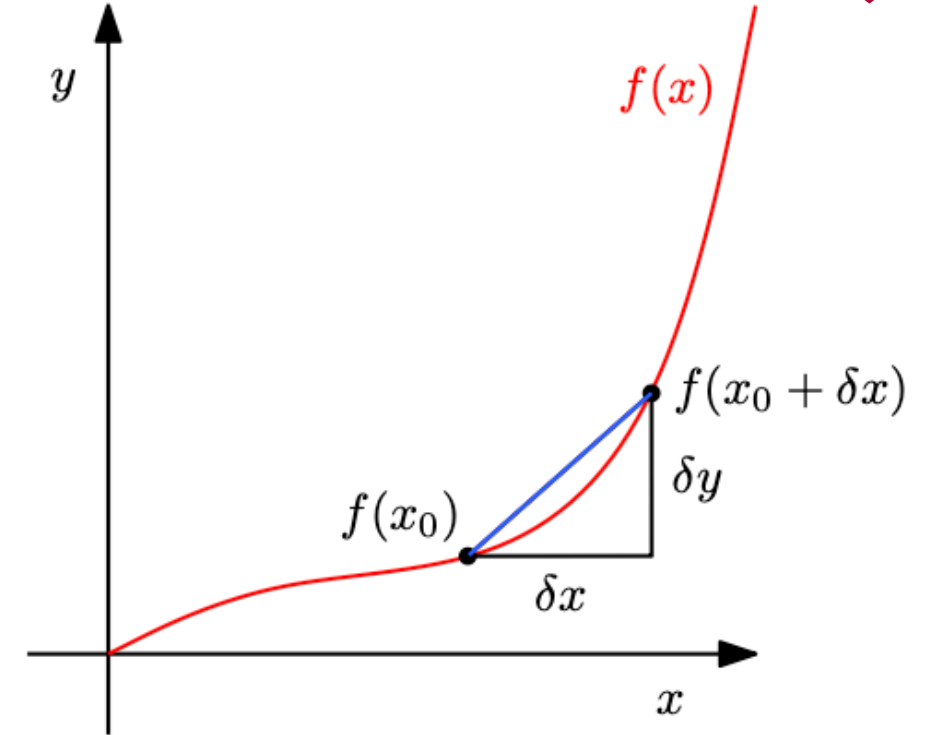
Derivative: The **derivative** of a function f at x is defined as the limit

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

and it points in the direction of steepest ascent of f .

Differentiation rules are

- Product rule: $(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$
- Quotient rule: $\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$
- Sum rule: $(f(x) + g(x))' = f'(x) + g'(x)$
- Chain rule:





2.3.2. Partial Differentiation and Gradients

Partial derivatives and Gradient: For a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^n$ of n variables x_1, \dots, x_n , the **partial derivatives** are defined as

$$\begin{aligned} \frac{\partial f}{\partial x_1} &= \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(\mathbf{x})}{h} \\ &\vdots \\ \frac{\partial f}{\partial x_n} &= \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(\mathbf{x})}{h} \end{aligned}$$

and collect them in the row vector

$$\nabla_{\mathbf{x}} f = \frac{df}{d\mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}$$



2.3.3. Matrix Derivatives

The derivative of a vector \mathbf{a} w.r.t. a scalar x is itself a vector whose components are given by

$$\left(\frac{\partial \mathbf{a}}{\partial x}\right)_i = \frac{\partial a_i}{\partial x}.$$

Similarly,

$$\left(\frac{\partial x}{\partial \mathbf{a}}\right)_i = \frac{\partial x}{\partial a_i} \text{ and } \left(\frac{\partial \mathbf{a}}{\partial \mathbf{b}}\right)_{ij} = \frac{\partial a_i}{\partial b_j}.$$

Furthermore,

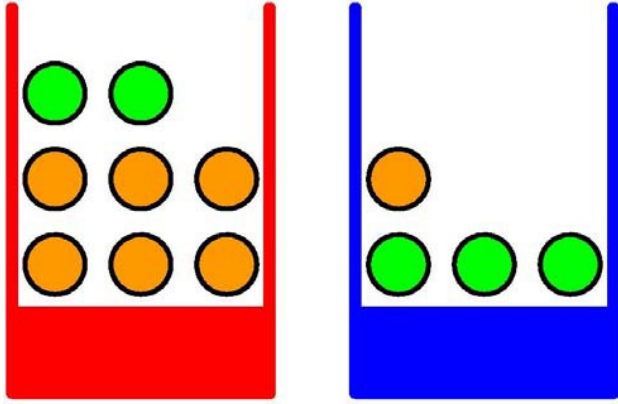
$$\begin{aligned}\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{a}) &= \frac{\partial}{\partial \mathbf{x}}(\mathbf{a}^T \mathbf{x}) = \mathbf{a} \\ \frac{\partial}{\partial \mathbf{x}}(\mathbf{A}\mathbf{B}) &= \frac{\partial \mathbf{A}}{\partial \mathbf{x}}\mathbf{B} + \mathbf{A}\frac{\partial \mathbf{B}}{\partial \mathbf{x}} \\ \frac{\partial}{\partial \mathbf{x}}(\mathbf{A}^{-1}) &= -\mathbf{A}^{-1}\frac{\partial \mathbf{A}}{\partial \mathbf{x}}\mathbf{A}^{-1}\end{aligned}$$



2.4. Probability Theory



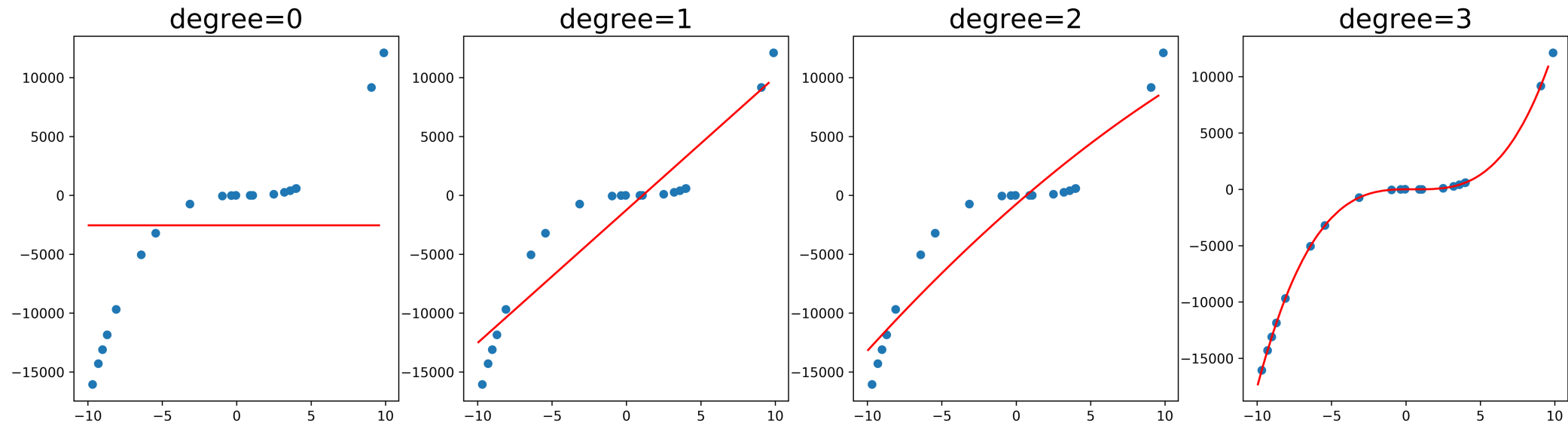
2.4.0. – Motivation: Probability Theory



- Let us look at the following example:
 - We have two boxes, one red and one blue
 - Red box: 2 apples and 6 oranges
 - Blue box: 3 apples and 1 orange
 - Pick red box 40% of the time and blue box 60% of the time, then pick one item of fruit
 - Question 1: what is the overall probability that the selection procedure will pick an apple?
 - Question 2: given that we have chosen an orange, what is the probability that the box we chose was the blue one?



2.4.0. Motivation: Curve Fit





2.4.2. Probability and Random Variables

- **Sample Space Ω :** the set of all possible outcomes of the experiment.
- **Event Space:** the space of potential results of the experiment. A subset of sample space is in the event space.
- **Probability:** measurements the probability or degree of belief that the event will occur.



2.4.2. Discrete Probabilities

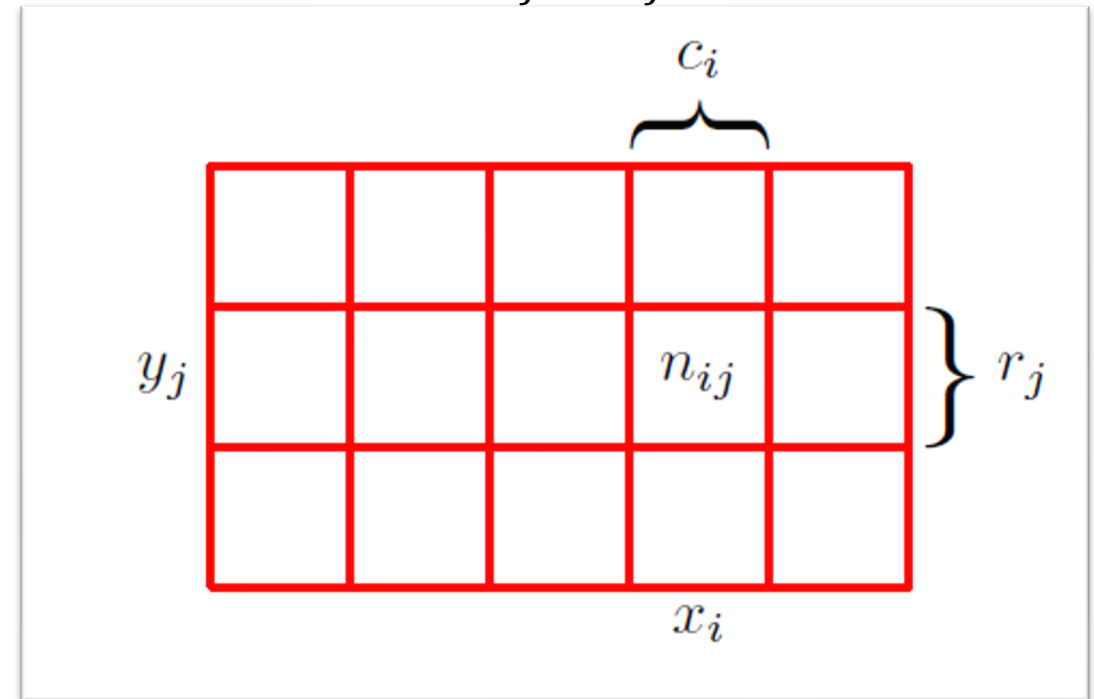
When the target space is *discrete*, the probability distribution of multiple random variables as filling out a (multidimensional) array of numbers.

Probability mass function: If x is a discrete variable, then $p(x)$ is sometimes called a *probability mass function* which implies as a set of probability masses concentrated at the allowed values of x .

2.4.2. Discrete Probabilities

Consider

- X can take any of the values x_i where $i = 1, \dots, M$.
- Y can take the values y_j where $j = 1, \dots, L$.
- A total of N trials in both of the variables X and Y .
- Let the number of trials for $X = x_i, Y = y_j$ be n_{ij} .
- Let the number of trials in which X takes the value x_i (irrespective of the value that Y takes) be denoted by c_i , and let the number of trials in which Y takes the value y_j be r_j .





2.4.2. Discrete Probabilities

- The probability that X will take the value of x_i and Y will take the value y_j , $p(X = x_i, Y = y_j)$, is called the *joint probability* of $X = x_i$ and $Y = y_j$.

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}.$$

- The probability that X takes the value x_i irrespective of the value of Y is

$$p(X = x_i) = \frac{c_i}{N}.$$

- **Sum Rule:** We have $c_i = \sum_j n_{ij}$ and therefore,

$$p(X = x_i) = \sum_{j=1}^L p(X = x_i, Y = y_j).$$

- Similarly, the probability that Y takes the value y_j irrespective of the value of X is

$$p(Y = y_j) = \frac{r_j}{N} = \sum_{i=1}^M p(X = x_i, Y = y_j).$$

- $p(X = x_i)$ and $p(Y = y_j)$ are sometimes called the *marginal* probability.



2.4.2. Discrete Probabilities

- Consider for which $X = x_i$, the fraction of instances for which $Y = y_j$ is $p(Y = y_j | X = x_i)$ called the *conditional* probability of $Y = y_j$ given $X = x_i$.
- It is obtained by finding the fraction of those points in i fall in (i, j) :

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}.$$

- The conditional probability of X given Y is

$$p(X = x_i | Y = y_j) = \frac{n_{ij}}{r_j}.$$



2.4.3. Continuous Probability

- Consider the target space with real continuous numbers \mathbb{R} .
- Most often, we pretend that we perform operations as we have discrete probability spaces with finite spaces. However, the simplification is not precise when
 - if operations were infinitely repeated. or
 - if a single point were drawn from an interval.

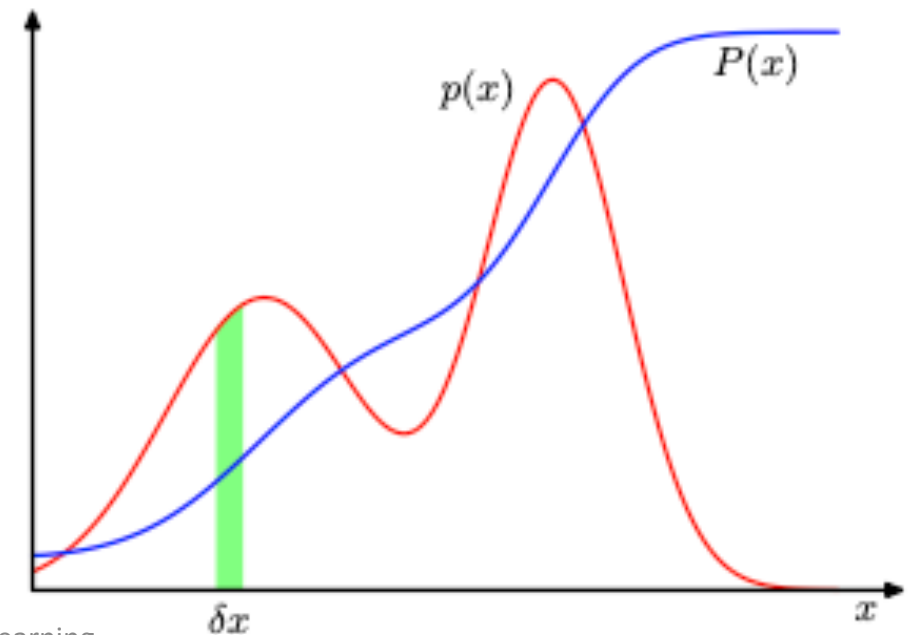
- **Probability Density Function** (pdf): If the probability of a real-valued variable x falling in the interval $(x, x + \delta x)$ is given by $p(x)\delta x$ for $\delta x \rightarrow 0$, then $p(x)$ is called the *probability density* over x .

- The probability that x will lie in an interval (a, b) is then given by

$$p(x \in (a, b)) = \int_a^b p(x) dx.$$

- The probability density $p(x)$ must satisfy the two conditions:

$$\begin{aligned} p(x) &\geq 0 \\ \int_{-\infty}^{\infty} p(x) dx &= 1. \end{aligned}$$





2.4.3. Continuous Probability

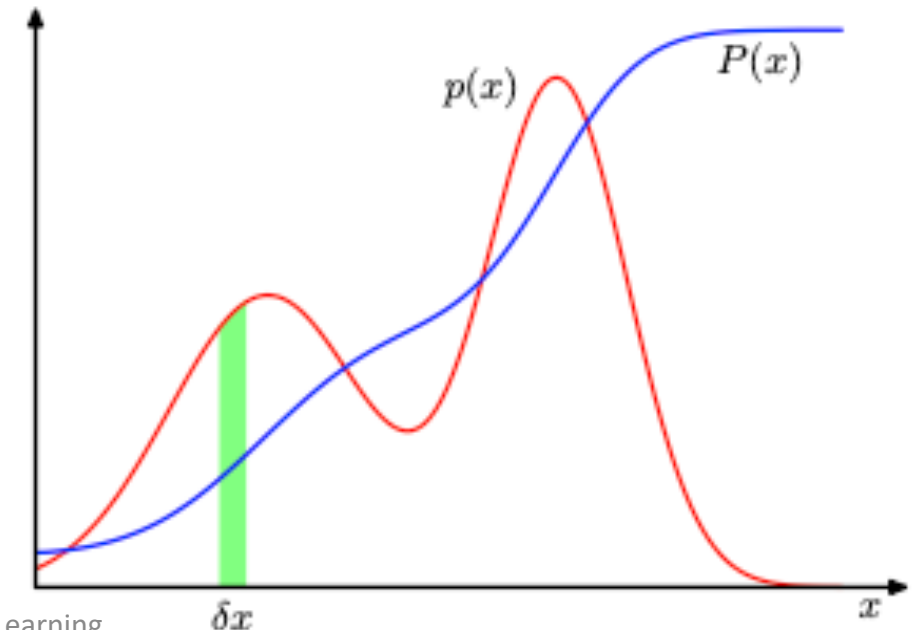
- If we consider a change variable $x = g(y)$, then a function $f(x)$ becomes $f(g(y))$. This means that a probability density $p_x(x)$ that corresponds to a density $p_y(y)$ with respect to the new variable y , where the suffixes denote the fact that $p_x(x)$ and $p_y(y)$ are different densities. The observation then transforms into the range $(y, y + \delta y)$ where $p_x(x)\delta x \sim p_y(y)\delta y$ and

$$p_y(y) = p_x(x) \left| \frac{dx}{dy} \right| = p_x(g(y)) |g'(y)|.$$

- **Cumulative distribution function (cdf):** The probability that x lies in the interval $(-\infty, z)$ is given by the *cumulative distribution function* defined by

$$P(z) = \int_{-\infty}^z p(x) dx$$

- Which satisfies $P'(z) = p(x)$.





2.4.4. Sum Rule, Product Rule, and Bayes' Theorem

Sum Rule: The sum rule is the *marginalization property*. It relates the joint distribution to a marginal distribution via

$$p(\mathbf{x}) = \begin{cases} \sum_{\mathbf{y} \in Y} p(\mathbf{x}, \mathbf{y}), & \text{if } \mathbf{y} \text{ is discrete} \\ \int_Y p(\mathbf{x}, \mathbf{y}) d\mathbf{y}, & \text{if } \mathbf{y} \text{ is continuous} \end{cases}.$$

Product Rule: It relates the joint distribution to the conditional distribution via

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x}).$$

- Every joint distribution of two random variables can be factorized of two other distributions.
- The ordering of random variables is arbitrary, $p(x, y) = p(x|y)p(y)$.



2.4.4. Sum Rule, Product Rule, and Bayes' Theorem

- In machine learning and Bayesian statistics, we are often interested in making inferences of unobserved random variables given that we have observed other random variables.
- Assume some prior knowledge $p(Y)$ about an observed random variable Y and some relationship $p(X|Y)$ between X and a second random variable Y , which we can observe.
- **Bayes' Theorem (rule or law):** From the product rule, together with the symmetry property $p(X, Y) = p(Y, X)$, we can build the following relationship called *Bayes' theorem* between conditional probabilities

$$\text{Posterior} = \text{Likelihood} \cdot \frac{\text{Prior}}{\text{Event}} \rightarrow p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}.$$

- We can view the denominator as the normalizer

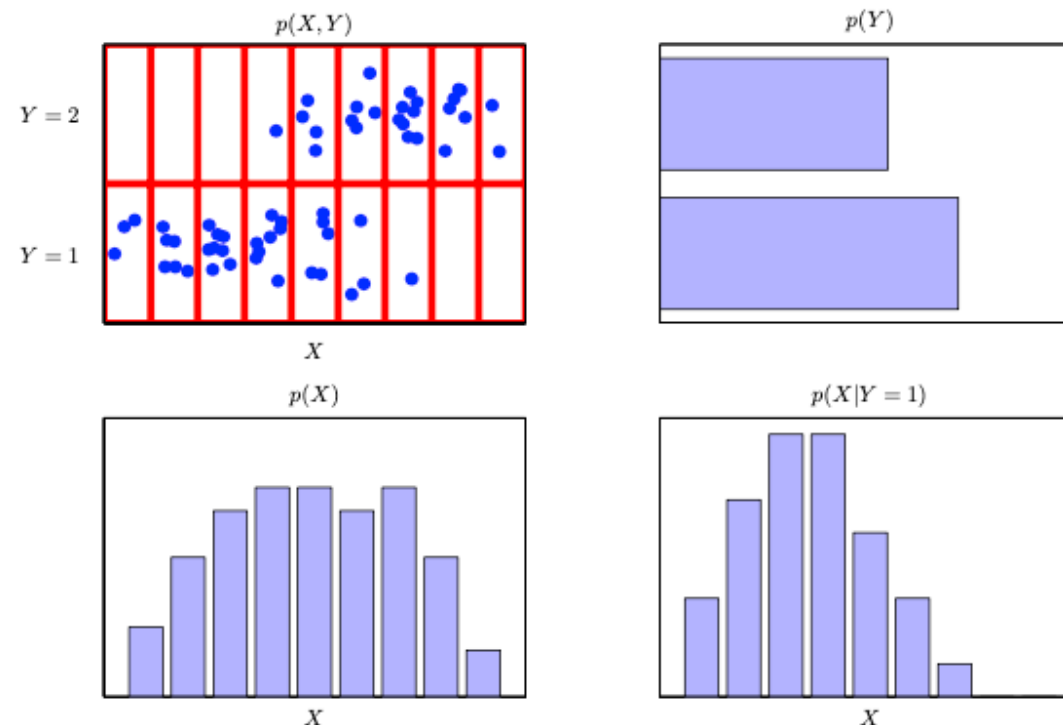
$$p(X) = \sum_y p(X|Y)p(Y)$$

using the sum rule.



2.4.4. Sum Rule, Product Rule, and Bayes' Theorem

- The *prior* $P(X)$ probability encapsulates the subjective prior knowledge of the unobserved variable X before observing any data.
- The *likelihood* $p(X|Y)$ describes how X and Y are related. It is the probability of the data X if the latent variable Y were to know. The likelihood is not a distribution in Y , but only in X .
- The *posterior* $P(Y|X)$ is the quantity of interest in Bayesian statistics.





2.4.4. Sum Rule, Product Rule, and Bayes' Theorem

- From a Bayesian perspective, probability provides a quantification of uncertainty.
- Considering the curve-fit example, we can use the probability theory to describe the uncertainty in model parameters and the models.
- Let the function be $t(x) = w_0 + w_1x + w_2x^2 + w_3x^3$.
 - Assume we know the prior probability of $\mathbf{w} = [w_0, w_1, w_2, w_3]$, $p(\mathbf{w})$.
 - The effect of the observed data $\mathcal{D} = \{t_1, \dots, t_N\}$ is expressed through the conditional probability $p(\mathcal{D}|\mathbf{w})$.
 - Then the uncertainty in \mathbf{w} after we have observed \mathcal{D} in the form of the posterior probability $p(\mathbf{w}|\mathcal{D})$:

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}.$$

- $p(\mathcal{D}|\mathbf{w})$ is a function of parameter vector \mathbf{w} called the *likelihood function* expressing how probable the observed data set is for different settings of the parameter vector \mathbf{w} .
- The denominator $p(\mathcal{D})$ is

$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) d\mathbf{w}.$$



2.4.4. Sum Rule, Product Rule, and Bayes' Theorem

- Given the definition of likelihood, we can state Bayes' theorem in words
$$\text{Posterior} \propto \text{likelihood} \times \text{prior}$$
where all of these quantities are viewed as functions of \mathbf{w} .
- In the sum and product rule approach, \mathbf{w} is considered to be a fixed parameter found by some 'estimator' and errors are obtained by considering the distribution of possible data sets \mathcal{D} .
- The common estimator is *maximum likelihood* in which \mathbf{w} is set to the value that maximizes $p(\mathcal{D}|\mathbf{w})$. In ML, we use the negative log of the likelihood function and call it the *error function* – it is a monotonically decreasing function. Maximizing the likelihood is equivalent to minimizing the error.
- In Bayesian approach, the uncertainty in the parameters is expressed through the probability distribution over \mathbf{w} in a single data set \mathcal{D} . Therefore, the inclusion of prior knowledge arises naturally.
- Example: Suppose that a fair-coin is tossed three times and lands heads each time. What is the probability of landing heads?



2.4.4. Sum Rule, Product Rule, and Bayes' Theorem

Question 1: What is the probability of picking an apple?

$$\begin{aligned}
 p(F = a) &= p(F = a|B = r)p(B = r) + p(F = a|B = b)p(B = b) \\
 &= \frac{2}{8} \left(\frac{4}{10} \right) + \frac{3}{4} \left(\frac{6}{10} \right) = \frac{11}{20}. \\
 p(F = o) &= 1 - p(F = a) = 1 - \frac{11}{20} = \frac{9}{20}.
 \end{aligned}$$

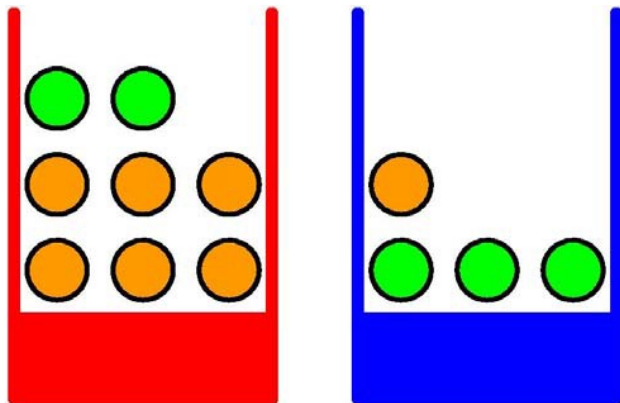
Question 2: What is the probability of picking an orange from the blue box?

$$\begin{aligned}
 p(B = b|F = o) &= \frac{p(F = o|B = b)p(B = b)}{p(F = o)} = 1 - p(r|o) = 1 - \frac{p(F = o|B = r)p(B = r)}{p(F = o)} \\
 &= 1 - \left(\frac{3}{4} \left(\frac{4}{10} \right) \left(\frac{20}{9} \right) \right) = \frac{1}{3}
 \end{aligned}$$

3/4

4/10

9/20





2.4.4. Sum Rule, Product Rule, and Bayes' Theorem

- **Prior and Posterior Probability:** The probability available before we observe the identity is called the prior probability and the probability after the observation is called the posterior probability. In this example, $p(B)$ is the prior probability and $p(B|F)$ is the posterior probability.
- **Independent:** If the joint distribution of two variables factorizes into the product of the marginals, so that

$$p(X, Y) = p(X)p(Y)$$

then, X and Y are *independent*. If each box had the same fraction of apples and oranges, then

$$p(F|B) = p(F).$$



2.4.5. Expectations and Covariances

- The average value of some function $f(x)$ under a probability distribution $p(x)$ is called the *expectation* of $f(x)$ and will be denoted as $\mathbb{E}(f)$

$$\mathbb{E}(f) = \sum_x p(x)f(x)$$

- So that the average is weighted by the relative probabilities of the different values of x . In the case of continuous variables, expectations are expressed in terms of an integration w.r.t. the corresponding probability density

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

- For a case of several variables,

$$\mathbb{E}_x(f(x, y))$$

where the subscript x is the variable being averaged over and $f(x, y)$ is the function w.r.t. the distribution of x .

- The *conditional expectation* w.r.t. a conditional distribution is

$$\mathbb{E}_x(f|y) = \sum_x p(x|y)f(x).$$



2.4.5. Expectations and Covariances

- The *variance* of $f(x)$ is defined by

$$\text{var}(f) = \mathbb{E} \left[\left(f(x) - \mathbb{E}(f(x)) \right)^2 \right]$$

- And provides a measure of how much variability there is in $f(x)$ around its $\mathbb{E}(f(x))$. The variance also can be written as $\mathbb{E}(f(x)^2) - \mathbb{E}(f(x))^2$:

$$\begin{aligned} \mathbb{E} \left[\left(f(x) - \mathbb{E}(f(x)) \right)^2 \right] &= \mathbb{E} \left[f(x)^2 - 2f(x)\mathbb{E}(f(x)) + \mathbb{E}(f(x))^2 \right] \\ &= \mathbb{E}[f(x)^2] - 2\mathbb{E}(f(x))\mathbb{E}(f(x)) + \mathbb{E} \left[\mathbb{E}(f(x))^2 \right] \\ &= \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2. \end{aligned}$$

- For two random variables x and y , the *covariance* is the product of their deviations from their respective means,

$$\text{cov}[x, y] = \mathbb{E}_{x,y}[\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] = \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y]$$

when x and y vary together. Otherwise, the covariance vanishes. We also can consider $\text{cov}(x) = \text{cov}[x, x]$.



2.4.5. Expectations and Covariances

- The *correlation* between two random variables is the normalized covariance between them via

$$\text{corr}(x, y) = \frac{\text{cov}[x, y]}{\sqrt{\text{var}(x)\text{var}(y)}}.$$

- The *empirical mean* vector is the arithmetic average of the observation for each variable and is defined as

$$\boldsymbol{\mu} = \bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

- The *empirical covariance* is

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T.$$

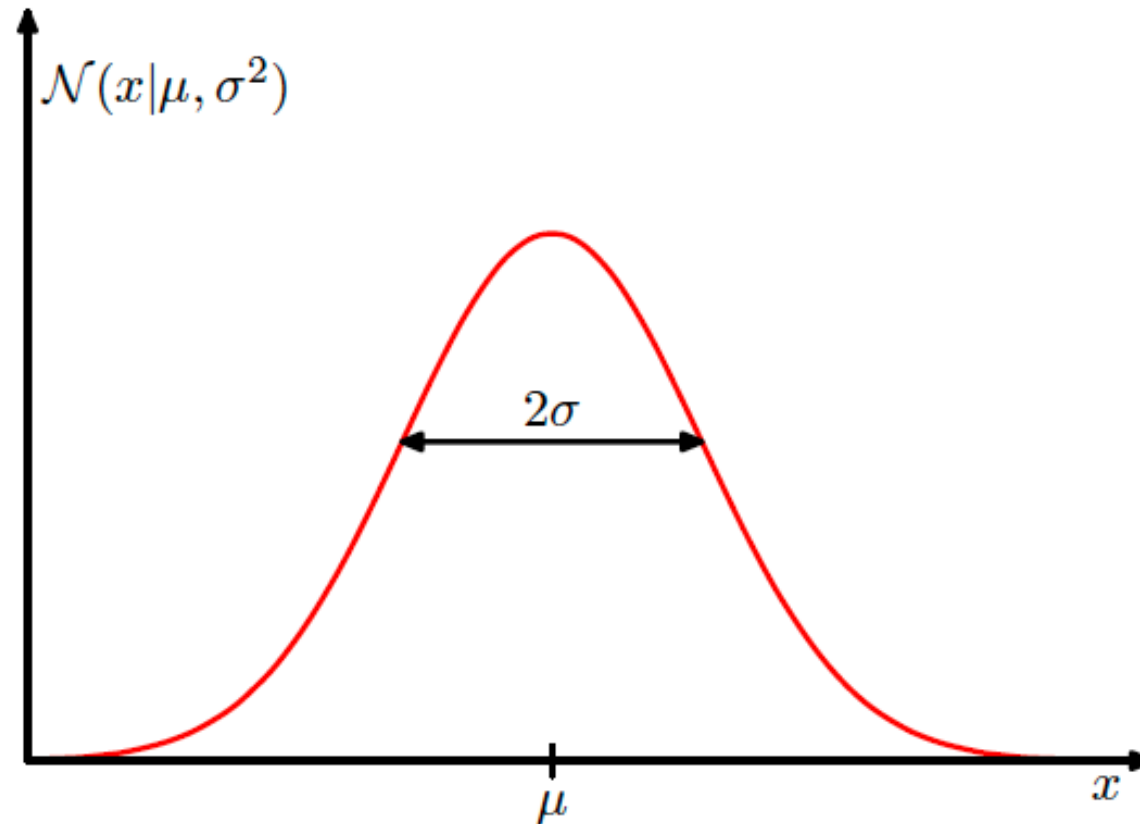


2.4.6. The Gaussian Distribution

For the case of a single real-valued variable x , the Gaussian distribution is defined by

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

where μ is the mean and σ^2 is the variance ($\sqrt{\sigma^2}$ is the standard deviation). Sometimes, we use $\beta = 1/\sigma^2$ and it is called the precision.





2.4.6. The Gaussian Distribution

Properties:

- Condition

$$\mathcal{N}(x|\mu, \sigma^2) > 0$$

- Normalization

$$\int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) dx = 1$$

- Expectation value

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x dx = \mu$$

- Variance

$$\text{var}[x] = \sigma^2$$

- The maximum of a distribution is the mode that coincides with the mean.



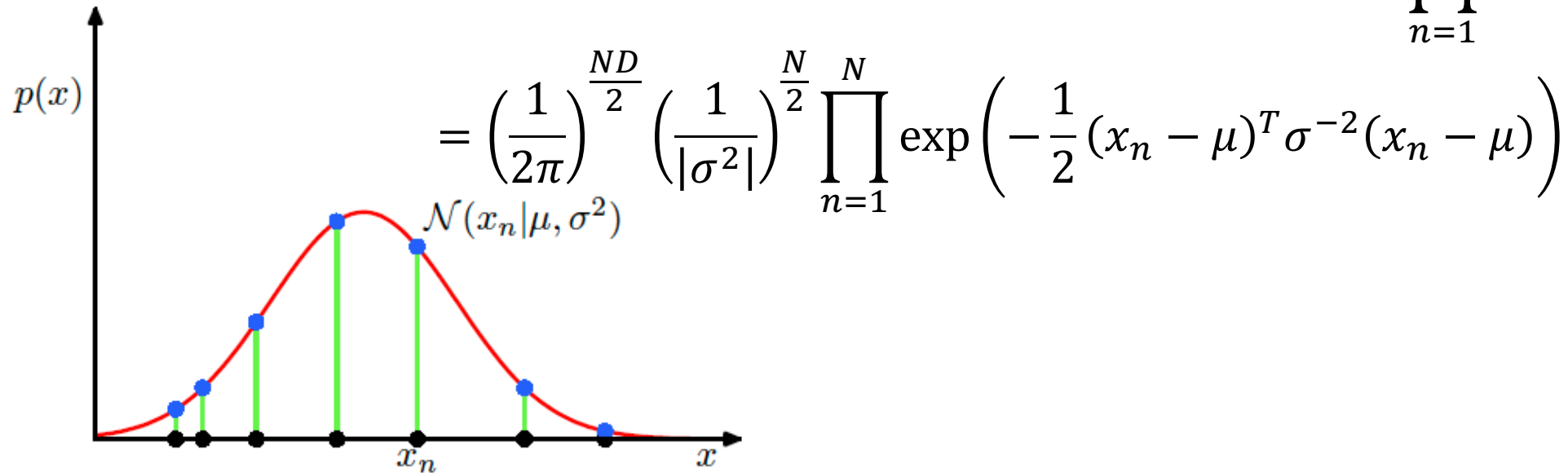
2.4.6. The Gaussian Distribution

- The Gaussian distribution over a D -dimensional vector \mathbf{x} of continuous variables is given by

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

- Suppose we have a data set of observations $X = (x_1, \dots, x_N)^T$ and the observations are drawn independently from a Gaussian distribution from unknown μ and σ^2 . We call this situation *independent and identically distributed* (i.i.d.). The probability of the data set is then

$$p(X|\mu, \sigma^2) = \mathcal{N}(x_1|\mu, \sigma^2)\mathcal{N}(x_2|\mu, \sigma^2) \cdots \mathcal{N}(x_N|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2)$$





2.4.6. The Gaussian Distribution

We can determine values for the unknown parameters by **maximizing the log-likelihood function**.

$$\ln p(X|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi).$$

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$$

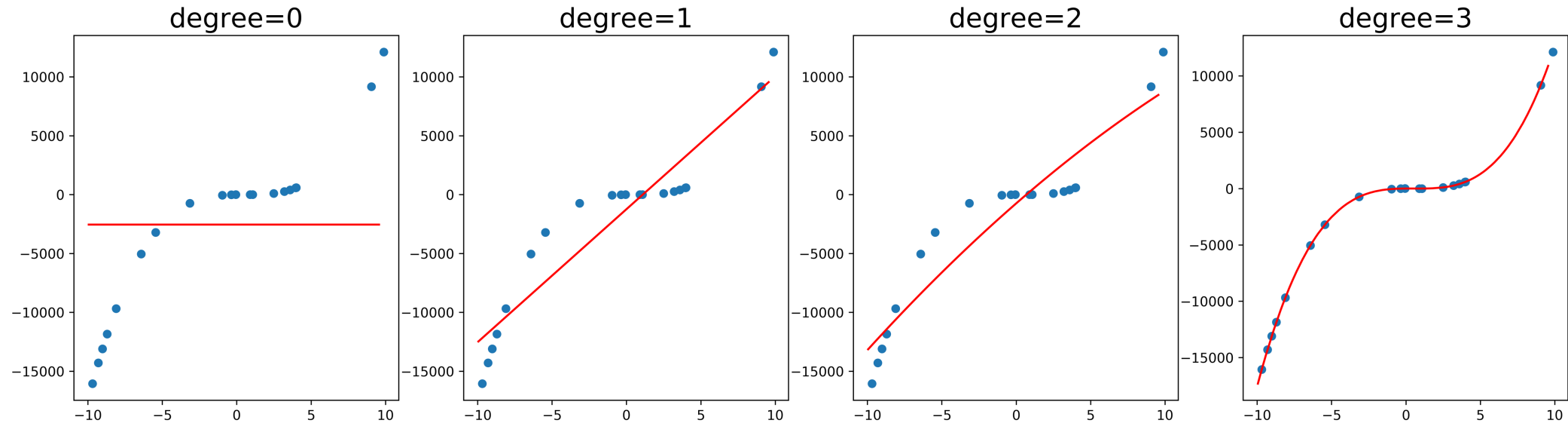
$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

which are the *sample mean* and *sample variance*, respectively. The maximum likelihood approach systematically underestimates the variance of the distribution.

The expectation of these parameters is

$$\begin{aligned} \mathbb{E}[\mu_{ML}] &= \mu \\ \mathbb{E}[\sigma_{ML}^2] &= \frac{N-1}{N} \sigma^2. \end{aligned}$$

2.4.7. Curve Fitting Revisit



- Suppose we predict for the target variable t given some new value of the input variable x on the basis of a set of data, $X = (x_1, \dots, x_N)^T$ and $T = (t_1, \dots, t_N)^T$. Assume that the given value of x , the corresponding value of t has a Gaussian distribution with a mean equal to the value $y(x, \mathbf{w})$ of the polynomial curve.

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1})$$



2.4.7. Curve Fitting Revisit

If the data are assumed to be i.i.d., then the likelihood function is

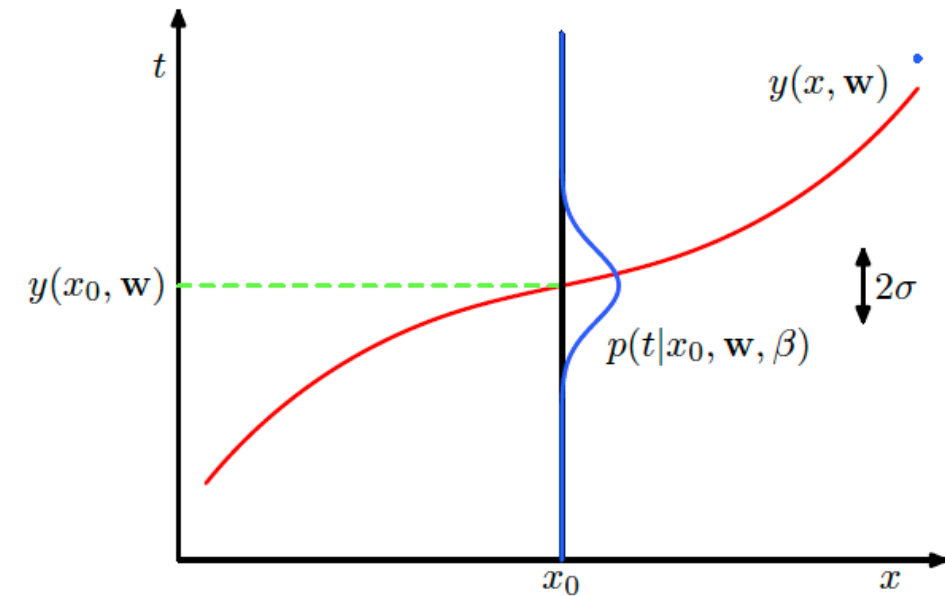
$$p(T|X, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1})$$

and the log likelihood function is

$$\begin{aligned} \ln p(T|X, \mathbf{w}, \beta) \\ = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi. \end{aligned}$$

If \mathbf{w} and β are determined, then we have a probabilistic model that are expressed in terms of the predictive distribution that gives the probability distribution over t , rather than simply a point estimate, and is obtained by substituting the maximum likelihood parameters

$$p(t|x, \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t | y(x, \mathbf{w}_{ML}), \beta_{ML}^{-1}).$$





2.4.7. Curve Fitting Revisit

If we use Bayesian approach with a prior distribution over the polynomial coefficients \mathbf{w} ,

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{\frac{M+1}{2}} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

where α is the precision of the distribution (called *hyperparameter* that controls the distribution of model parameter), and $M + 1$ is the total number of elements in the vector \mathbf{w} for an M^{th} order polynomial. The posterior distribution for \mathbf{w} is

$$p(\mathbf{w}|X, T, \alpha, \beta) \propto p(T|X, \mathbf{w}, \beta)p(\mathbf{w}|\alpha).$$

We can find \mathbf{w} by maximizing the posterior distribution called *maximum posterior* (MAP) by taking the negative log. The maximum of the posterior is equivalent to the minimum of the negative log of the posterior as shown

$$-\ln p(\mathbf{w}|X, T, \alpha, \beta) = -\ln p(T|X, \mathbf{w}, \beta) - \ln p(\mathbf{w}|\alpha).$$



2.4.8. Binary Variables

Consider a single binary random variable $x \in \{0,1\}$. The probability of $x = 1$ will be denoted by the parameter μ so that

$$p(x = 1|\mu) = \mu$$

where $0 \leq \mu \leq 1$, from which it follows that $p(x = 0|\mu) = 1 - \mu$. The probability distribution (*Bernoulli* distribution) over x can therefore be written in the form

$$\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x}.$$

This distribution is normalized and has mean and variance given by

$$\begin{aligned}\mathbb{E}[x] &= \mu \\ \text{var}[x] &= \mu(1 - \mu).\end{aligned}$$

Suppose we have a data set $D = \{x_1, \dots, x_N\}$. Assume that observations are i.i.d. from $p(x|\mu)$, so that

$$p(D|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n}(1 - \mu)^{1-x_n}.$$

The log likelihood is

$$\ln p(D|\mu) = \sum_{n=1}^N \ln p(x_n|\mu) = \sum_{n=1}^N \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\}.$$



2.4.9. Multinomial Variable

- Consider the Bernoulli distribution to an K -dimensional binary variable $x_k \in \{0,1\}$ such that $\sum_k x_k = 1$. Then the distribution of \mathbf{x} is given

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k}$$

where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T$, and the parameters μ_k are constrained to satisfy $\mu_k \geq 0$ and $\sum_k \mu_k = 1$.



2.5. Conclusion



2.5.2. Conclusion

- A brief review of Mathematics that will play central roles in Machine Learning Algorithms
 - Linear Algebra, Analytic Geometry, and Vector Calculus will help to understand the work of algorithms.
 - Probability theory will help to understand the characteristics of algorithms
 - If discussed topics are not fully digestible, it is still okay.
 - Most of the terms will be repeatedly discussed throughout the semesters.