

STATISTICAL AND MATHEMATICAL METHODS FOR AI:

ERRORS: discrete representations for continuous numbers lead to approximation errors (measurement, arithmetic, truncation, inherent)

ABSOLUTE ERROR: let x be a ground truth and \tilde{x} be the result of some operation, then $E_x := |\tilde{x} - x|$

RELATIVE ERROR: $R_x := \frac{E_x}{|x|}$, facilitates understanding the error's order of magnitude

ROUNDING RULES: rules for approximation of $x \in \mathbb{R}$ to $f(x) \in \mathbb{F}$

Truncation: $f_l(x) = x_p$ truncated after the p -th digit

Rounding: $f_l(x) = y \in \mathbb{F}$ nearest to $x \in \mathbb{R}$, most accurate default IEEE

FLOATING POINT SYSTEM: a floating point system is defined as $\mathbb{F} = (\beta, I, L, U)$ where β = base, I = precision, and $[L, U]$ = exponent range

NORMALIZED SCIENTIFIC NOTATION: let $x \in \mathbb{R}$ s.t. $a_1, \dots, a_n \neq 0$ are the first n non-zero digits of x . Then $f_l(x) = 0.a_1 \dots a_n \cdot \beta^k$ is the normalized representation $\mathbb{F}(\beta, n, L, U)$ for x

MACHINE PRECISION: accuracy of the floating point system, $\min(\varepsilon)$ s.t. $f_l(1 + \varepsilon) > 1$, maximum relative representation error ($\forall x \in \mathbb{R}$) $\frac{|f_l(x) - x|}{|x|}$)

N.B.: 1) $\varepsilon = \beta^{-t}$ with rounding by truncation

2) Arithmetic operation $\odot: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, $\odot: \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$ causes a rounding error corresponding to $\left| \frac{(x \odot) - (x \cdot y)}{(x \cdot y)} \right| < \varepsilon$

LINEAR SYSTEMS: $A\bar{x} = b \Leftrightarrow \exists A' (A \text{ non singular}) \Leftrightarrow \det(A) \neq 0 \Leftrightarrow x = A'^{-1}b$, A' is an expensive computation $O(n^3)$

LU FACTORIZATION: direct method for linear systems $A\bar{x} = b$, let $A = LU$ (lower and upper triangulars) Then $A\bar{x} = b \equiv L\bar{U}\bar{x} = b$ thus $L\bar{y} = b$ and $U\bar{x} = \bar{y}$

N.B.: unstable, subject to large approximation errors in computing L and U .

PIVOTING: let $A' = PA$ where $P = I$ with swapped rows Then $A\bar{x} = b \equiv A'\bar{x} = P\bar{b}$ Thus $L\bar{U}\bar{x} = P\bar{b}$, $L\bar{y} = P\bar{b}$ and $U\bar{x} = \bar{y}$

CHOLESKY FACTORIZATION: let $A \in \mathbb{R}^{n \times n}$ s.t. $A = L^T L$ and $A\bar{x} = b \equiv L\bar{L}^T\bar{x} = b$ Thus $L\bar{y} = b$ and $L\bar{x} = \bar{y}$

AFFINE SPACES: let V be a vector space, $x_0 \in V$ and $U \subseteq V$ be a subspace Then $L := x_0 + U = \{x_0 + u \mid u \in U\} = \{v \in V \mid \exists u \in U \text{ s.t. } v = x_0 + u\} \subseteq V$ is said to be an

affine subspace (linear manifold) of V , furthermore U and x_0 are said to be direction and support point of L

N.B.: 1) if $x_0 \notin U$ then $\bar{0} \notin L$ Thus $L \neq \text{vector subspace of } V$ for $x_0 \notin U$.

2) affine subspaces can be described by parameters, indeed $\forall x \in L \quad x = x_0 + \lambda_1 b_1 + \dots + \lambda_n b_n$ where (b_1, \dots, b_n) = ordered basis of U .

AFFINE MAPPINGS: let $\Phi: V \rightarrow W$ be a linear mapping and $w \in W$, then $\Psi: V \rightarrow W$ s.t. $\Psi(x) = w + \Phi(x)$ is said to be an affine mapping

N.B.: w is said to be the translation vector of Ψ .

NORM: $\|\cdot\|: V \rightarrow \mathbb{R}$ assigns a notion of length s.t.:

$$1) \forall x \in V \quad \|x\| \geq 0 \wedge \|x\| = 0 \Leftrightarrow x = 0$$

$$2) \|\lambda \cdot x\| = |\lambda| \cdot \|x\|$$

$$3) \|x + y\| \leq \|x\| + \|y\|$$

MANHATTAN NORM: $\|x\|_1 := \sum_{i=1}^n |x_i|$

EUCLIDEAN NORM: $\|x\|_2 := \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{x^T x}$

N.B.: in general $\|\cdot\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$

DOT PRODUCT: let $x, y \in V$, the dot product of x and y is defined as $\langle x, y \rangle := \sum_{i=1}^n x_i y_i$

MATRIX NORM: The norm of a matrix is defined as a function $\|\cdot\|: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ s.t.

$$1) \forall A \in \mathbb{R}^{n \times m} \quad \|A\| \geq 0 \wedge \|A\|=0 \Leftrightarrow A=0$$

$$2) \forall \alpha \in \mathbb{R}, A \in \mathbb{R}^{n \times m} \quad \|\alpha A\| = |\alpha| \|A\|$$

$$3) \forall A, B \in \mathbb{R}^{n \times m} \quad \|A+B\| \leq \|A\| + \|B\|$$

FROBENIUS NORM: $\|A\|_F := \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$

2-NORM: $\|A\|_2 := \sqrt{\rho(A^T A)}$ where $\rho(A) := \max(\text{eigvals}(A))$ (= spectral radius)

POSITIVE DEFINITE MATRIX: let $A \in \mathbb{R}^{n \times n}$ we call A positive definite if $\forall x \in V \setminus \{0\} \quad x^T A x > 0$, A is positive semi-definite if $x^T A x \geq 0$

EIGENVALUES: let $A \in \mathbb{R}^{n \times n}$ and $\lambda \in \mathbb{C}$ be an eigenvalue of A then x is the corresponding eigenvector if $Ax = \lambda x$

N.B.: 1) let $A \in \mathbb{R}^{n \times n}$ symmetric $\wedge \text{eigvals}(A) > 0$ then $A > 0$, if $\text{eigvals}(A) \geq 0$ then $A \geq 0$

2) let $A \in \mathbb{R}^{m \times n}$ then $A^T A \geq 0 \in \mathbb{R}^{n \times n}$ symmetric, $A^T A \geq 0 \Leftrightarrow \text{rank}(A^T A) = \min(m, n)$ (full rank)

3) let $A \in \mathbb{R}^{n \times n}$ then $\text{rank}(A) = n \Leftrightarrow \forall \lambda \in \text{eigvals}(A), \lambda \neq 0$

ORTHOGONAL VECTORS: let $x, y \in \mathbb{R}^n$, x and y are said to be orthogonal if $\langle x, y \rangle = 0$ ($x \perp y$)

ORTHONORMAL VECTORS: let $x, y \in \mathbb{R}^n$, x and y are said to be orthonormal if $\langle x, y \rangle = 0 \wedge \|x\| = \|y\| = 1$

ORTHOGONAL MATRICES: $A \in \mathbb{R}^{n \times n}$ is said to be orthogonal if its columns are orthonormal

N.B.: 1) let A be orthogonal then $A A^T = I = A^T A$, i.e. $A^T = A^{-1}$

2) Transformations by orthogonal matrices don't change a vector's length

ORTHONORMAL BASIS: let $V \subset \mathbb{R}^n$, $\{b_1, \dots, b_m\}$ be a basis of V , if $\forall i, j \in \{1, \dots, n\} \quad \langle b_i, b_j \rangle = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$ then $\{b_1, \dots, b_m\}$ is called an orthonormal basis.

ORTHOGONAL COMPLEMENT: let V s.t. $\dim(V) = D$ and $U \subset V$ s.t. $\dim(U) = M$, we call orthogonal complement of U $U^\perp := \{x \in V \mid \exists y \in U \text{ s.t. } x \perp y\}$ s.t. $\dim(U^\perp) = D - M$

N.B.: $U \cap U^\perp = \{0\}$ i.e. $\forall x \in V \quad x = \sum_{i=1}^M \lambda_i b_i + \sum_{j=1}^{D-M} \psi_j b_j^*$ where $\lambda_i, \psi_j \in \mathbb{R}$, (b_1, \dots, b_n) basis of U and (b_1^*, \dots, b_{D-M}) basis of U^\perp , therefore let w be the basis of U^\perp we have $w \perp U$ (normal vector)

PROJECTIONS: let V be a vector space and $U \subset V$ a linear mapping $\pi: V \rightarrow U$ is called a projection if $\pi^2 = \pi$ ($\pi^2 = \pi$ on U)

PROJECTION MATRIX: a projection matrix is defined as P_π s.t. $P_\pi^2 = P_\pi$ (= projection expressed as a matrix)

N.B.: 1) let $x \in \mathbb{R}^n$ then $\pi_U(x) = u \in U$ s.t. $\|x - u\|$ is minimal, let $Ax = b$ s.t. $Ax \neq b$ then $\pi_U(b)$ can be used to find an approximated solution

2) $\pi_U(x) \in U$, thus $\pi_U(x) = \lambda b$ where $\lambda \in \mathbb{R}$ and b is the basis vector of U .

3) let $U \subset \mathbb{R}^n$ s.t. $\dim(U) = m \geq 1$ then $\pi(x) = \sum_{i=1}^m \lambda_i b_i$

ALGEBRAIC MULTIPLICITY: let $\lambda_i \in \text{eigvals}(A)$ the algebraic multiplicity of λ_i is the number of times $\sqrt[n]{\lambda_i}$ appears in the characteristic polynomial of A

THEOREM: let $A \in \mathbb{R}^{n \times n}$, $\{\lambda_1, \dots, \lambda_n\} = \text{eigvals}(A)$ s.t. $\lambda_i \neq \lambda_j \forall i \neq j$ and x_1, \dots, x_n be the corresponding eigenvectors. Then x_1, \dots, x_n are linearly independent.

DEFECTIVE MATRIX: let $A \in \mathbb{R}^{n \times n}$ we call A defective if it has less than n linearly independent eigenvectors.

N.B.: a non-defective matrix may have less than n distinct eigenvalues.

DIAGONALIZABLE MATRIX: let $A \in \mathbb{R}^{n \times n}$ is diagonalizable if $\exists P \in \mathbb{R}^{n \times n}$ s.t. $D = P^{-1}AP$ is diagonal, A is similar to D .

THEOREM: let $A \in \mathbb{R}^{n \times n}$ Then A can be factored into $A = P^{-1}DP$ s.t. $P \in \mathbb{R}^{n \times n}$ and $D \in \mathbb{R}^{n \times n}$ diagonal s.t. diag(D) = eigenvals(A) \Leftrightarrow The eigenvectors of A are a basis of \mathbb{R}^n

SPECTRAL THEOREM: let $A \in \mathbb{R}^{n \times n}$ symmetric Then \exists b orthonormal basis of $V \subseteq \mathbb{R}^n$ s.t. b_1, \dots, b_n are eigenvectors of A

SINGULAR VALUE DECOMPOSITION (SVD): let $A \in \mathbb{R}^{m \times n}$ rectangular $\wedge \text{rank}(A) = r \in [0, \min(m, n)]$ Then $\exists U \in \mathbb{R}^{m \times m}$ orthogonal, $V \in \mathbb{R}^{n \times n}$ orthogonal and $\Sigma \in \mathbb{R}^{m \times n}$ where

$\Sigma_{ii} = \sigma_i \geq 0 \wedge \Sigma_{ij} = 0 \forall i \neq j$ s.t. $A = U\Sigma V^T$; $\sigma_1, \dots, \sigma_r$ are called singular values of A , by convention: $\sigma_1 \geq \dots \geq \sigma_r \geq 0$

N.B.: 1) cols(U) and cols(V) are called respectively left and right singular vectors of A , orthonormal bases of \mathbb{R}^m and \mathbb{R}^n .

2) We know $\forall A \in \mathbb{R}^{m \times n}$ rectangular $\rightarrow A^T A \geq 0$ symmetric, furthermore by the spectral theorem $A^T A$ is diagonalizable, thus $\exists P$ orthogonal s.t. $A^T A = P D P^T$ where $D = \text{diag}(\lambda_1, \dots, \lambda_n)$

3) $\exists \text{SVD}(A) \rightarrow A^T A = (U\Sigma V^T)(U\Sigma V^T)^T = V\Sigma^T U^T U \Sigma V^T = V\Sigma^T I \Sigma V^T = V\Sigma^2 V^T$, thus $A^T A$ is similar to Σ^2 and $\sigma_i^2 = \lambda_i$ where $\lambda_i \in \text{eigenvals}(A^T A)$ and $\sigma_i = \sqrt{\lambda_i}$

4) $\sigma_i = \sqrt{\lambda_{\max}} = \sqrt{\lambda_i(A^T A)} = \|A\|_F$, $\sigma_i = \sqrt{\lambda_i}$ i.e. $\|A\| = 1/\sigma_r(A)$

MATRIX APPROXIMATION: let $\text{rank}(A) = r$ and $A = U\Sigma V^T$ then $\hat{A}_k = \text{outer}(u_k v_k^T)$ where u_k, v_k = columns of U and V , $\text{rank}(\hat{A}_k) = k$, $A = \sum_{i=1}^r \sigma_i u_i v_i^T = \sum_{i=1}^r \sigma_i A_i$, Then $\hat{A}_k(k) = \sum_{i=1}^k \sigma_i A_i$ with $k \leq r$ is a k -rank approximation of A

N.B.: a higher rank matrix gives a better representation of the space it belongs to (rank = number of vectors forming a basis)

THEOREM: let $A, B \in \mathbb{R}^{m \times n}$ s.t. $\text{rank}(A) = r$ and $\text{rank}(B) = k$ Then $\forall k \leq r$ $\hat{A}_k(k) = \hat{A}_k(k) = \sum_{i=1}^k \sigma_i A_i = \arg \min_B \|A - B\|_F$, moreover $\|A - \hat{A}_k(k)\|_F = \sigma_{k+1}$

N.B.: if $\sigma_1, \dots, \sigma_n$ are known then k should be chosen as $\text{argmax}(\sigma_k - \sigma_{k+1})$ to minimize the relative loss of information in the approximated values

PRINCIPAL COMPONENT ANALYSIS (PCA): dimensionality reduction technique, let $X \in \mathbb{R}^{d \times N}$ be a high-dimensional dataset, define $Z \in \mathbb{R}^{k \times N}$ where $k \leq d$ s.t. $Z = P X$ where $P = U_k \in \mathbb{R}^{k \times d}$

N.B.: the columns of Z are said to be the principal components of X .

BACKPROPAGATION: minimize loss function of a model through gradient descent, given a loss function L , compute ∇L w.r.t. the model's parameters:

let $y = (f_k \circ \dots \circ f_1)(x) = f_k(\dots(f_1(x))\dots)$ be the true result of a multi-level function and $f_i(x_{i-1}) = \sigma(A_{i-1}x_{i-1} + b_{i-1})$ s.t. σ is an activation function and x_{i-1} be the output of the i -th layer, then L is a measure of the distance between y and $f_k(\theta, x)$ where $\theta = \{A_{0,0}, \dots, A_{k-1, k-1}\}$ (=parameters), Thus $\nabla L = \frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial f_k} \frac{\partial f_k}{\partial f_{k-1}} \dots \frac{\partial f_1}{\partial \theta}$, $\forall i$

AUTOMATIC DIFFERENTIATION: generalization of backpropagation, set of techniques to numerically evaluate the exact gradient of a function using intermediate variables obtained through

a sequence of elementary operations resulting in automatic computation of the gradient even for complex functions

N.B.: 1) backpropagation applies automatic differentiation in the opposite direction of the data flow.

2) not all programs can be automatically differentiated

OPTIMIZATION: find optimum of a function, machine learning models can be expressed as numerical optimization methods.

OPTIMALITY CONDITIONS:

First Order: let f be differentiable around x^* local minimum then $\nabla f(x^*) = 0$ (stationary point)

Second Order: 1) let f be twice differentiable around a local minimum x^* then $\nabla f(x^*) = 0$ and $Hf(x^*) \geq 0$

2) let f be twice differentiable around x^* s.t. x^* is stationary $\wedge Hf(x^*) > 0$ Then x^* is a local minimum

CONVEX SET: a set $C \subseteq \mathbb{R}^n$ is said to be convex if $\forall x, y \in C, \theta \in [0, 1], \theta x + (1-\theta)y \in C$ (\forall point p between x and y $p \in C$)

CONVEX FUNCTION: let C be convex and $f: C \rightarrow \mathbb{R}$, f is said to be convex if $\forall x, y \in C$ $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$ ($\forall p$ between $f(x)$ and $f(y)$ lies above the function)

N.B.: let f be convex Then \forall local minimum x , x is a global minimum. $\nabla f(x) = 0$ Then x = global minimum

LEVEL CURVES: $L_c := \{x \in \mathbb{R}^n | f(x) = c\}$

DESCENT METHODS: solve $\min(f(x))$ by generating $x_k = x_{k-1} + \alpha_k p_k$ where $p_k \in \mathbb{R}$ is called the search direction and $\alpha_k \in \mathbb{R}$ is called the step length.

Descent Direction: choose p_k s.t. $f(x_k) < f(x_{k-1})$, let f be continuous and differentiable in x and $p \in \mathbb{R}^n \neq 0$ s.t. $p^\top \nabla f(x) < 0$. Then p is a descent direction for f in x , $-\nabla f(x)$ is a descent direction.

N.B.: let p_k be a descent direction for f in x_k . Then $\forall \alpha_k \in [0, \bar{\alpha}]$ $f(x_k + \alpha_k p_k) < f(x_{k-1})$ for some $\bar{\alpha}$.

Step Length:

Constant: doesn't guarantee convergence

Inexact Search: choose α_k to produce sufficient decrease or a sufficient shift from $f(x_{k-1})$

Backtracking: let $\bar{\alpha} = 1$ (starting value) and $c \in [0, 1]$ if $f(x_k + \bar{\alpha} p_k) < f(x_k) + c \bar{\alpha} p_k^\top \nabla f(x_k)$ then $\alpha_k = \bar{\alpha}$ else $\bar{\alpha} = \bar{\alpha} \cdot p$ (usually $p=1/2$)

Wolfe Condition: $f(x_k + \bar{\alpha} p_k) < f(x_k) + c \bar{\alpha} p_k^\top \nabla f(x_k)$ (sufficient decrease, Armijo Condition), $-p_k^\top \nabla f(x_k + \alpha_k p_k) \leq -c_2 p_k^\top \nabla f(x_k)$ (curvature condition) satisfied by backtracking

THEOREM: if α_k is chosen by backtracking. Then the descent method converges to a local minimum.

GRADIENT DESCENT: let f be continuous and differentiable given a starting point x_0 . $\forall k$ select $\alpha_k = c$ via backtracking and $p_k = -\nabla f(x_k)$, iterate until conditions are met

N.B.: ADAM is a gradient descent variant (different rules for α_k)

STOP CONDITIONS: $\|\nabla f(x_k)\| < \gamma_1$ (or $\frac{\|\nabla f(x_k)\|}{\|\nabla f(x_0)\|} < \gamma_1$) \wedge $\|x_k - x_{k-1}\| < \gamma_2$ (or $\frac{\|x_k - x_{k-1}\|}{\|x_{k-1}\|} < \gamma_2$) \vee iterations $> K$

STARTING GUESS: usually random, influences result

CONVERGENCE SPEED: if $\|x_k - x^*\| < \|x_k - x^*\|^Q \cdot c$ for some c , then the method has convergence speed $= Q$

N.B.: $Q > 0$, $Q=1 \rightarrow$ linear speed $\wedge c < 1$, $Q \in [1, 2] \rightarrow$ superlinear speed, $Q=2 \rightarrow$ quadratic speed.

GRADIENT DESCENT WITH MOMENTUM: $x_{k+1} = x_k - \alpha_k (\nabla f(x_k))^\top + \gamma \Delta x_k$ where $\Delta x_k = x_k - x_{k-1}$ and $\gamma \in [0, 1]$ (x_k depends on 2 previous steps)

STOCHASTIC GRADIENT DESCENT: commonly used in neural networks, in general given N samples, $L(\theta, X, Y) = \sum_{n=1}^N L_n(\theta, x_n, y_n)$ Thus $\nabla_\theta L(\theta, X, Y) = \sum_{n=1}^N \nabla_\theta L_n(\theta, x_n, y_n)$,

let $S \subseteq \{1, \dots, N\}$ approximate $\nabla_\theta L(\theta, X, Y) \approx \sum_{i \in S} \nabla_\theta L_i(\theta, x_i, y_i)$, choosing a different random subset at each iteration. Let $k = |S|$ Then S is said to be a minibatch of size k

Epoch: \neq iteration, an epoch ends every time the entire dataset has been used, stop condition = maximum number of epochs (\wedge convergence checks, \wedge backtracking)

DISCRETE DISTRIBUTIONS:

Probability Mass Function: let X be a discrete random variable, The PMF of X is a function $f_X(x): X \rightarrow [0, 1]$ s.t. $f_X(x) = P(X=x)$

N.B.: 1) a PMF is a probability

2) let X be discrete $E[X] = \sum_{x_i \in X} x_i f_X(x_i) (= \mu)$, $\sigma^2 = \sum_{x_i \in X} (x_i - \mu)^2 f_X(x_i)$

Uniform Distribution: let X be a discrete random variable s.t. $f_X(x_i) = \frac{1}{|X|} \quad \forall x_i \in X$ Then $X \sim \text{Unif}(a, b)$ where $a = \inf(X)$ and $b = \sup(X)$

Poisson Distribution: let X be a discrete random variable and $\lambda \in \mathbb{R}$ s.t. $f_X(x_i) = \frac{x_i e^{-\lambda}}{x_i!}$ Then $X \sim \text{Poisson}(\lambda)$

N.B.: 1) The parameter λ represents the expected number of events in a given interval, The PMF represents the probability of x_i events occurring in the given interval.

2) $E[X] = \sigma^2 = \lambda$

CONTINUOUS DISTRIBUTIONS:

Probability Density Function: let $X: \Omega \rightarrow [a, b]$ be a continuous random variable, The PDF of X is a function $f_X: X \rightarrow \mathbb{R}$ s.t $P(X=x_i \in [c, d]) = \int_c^d f_X(x_i) dx_i$

N.B.: $f_X(x_i) > 0 \forall x_i \in X, \int_X f_X(x_i) dx_i = 1$, PDF is not a probability, its integral is (#PMF)

2) let X be continuous $E[X] = \int_X x_i f_X(x_i) dx_i, \sigma^2 = \int_X (x - \mu)^2 f_X(x_i) dx_i, \sigma = \sqrt{\sigma^2}$ is called standard deviation

Uniform Distribution: let X be continuous in $[a, b]$ s.t. $f_X(x_i) = \frac{1}{b-a} \forall x_i \in X$ Then $X \sim U(a, b)$

Gaussian Distribution: let X be continuous s.t. $f_X(x_i) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ Then $X \sim N(\mu, \sigma)$

MULTIVARIATE DISTRIBUTIONS:

Joint Probability: let X, Y be discrete, then their joint probability is given by $f_{XY}: X \times Y \rightarrow [0, 1]$ s.t. $f_{XY}(x_i, y_i) = P(X=x_i, Y=y_i)$

Marginal Probability: let (X, Y) be a random vector, its marginal probabilities are given by $p_X(x_i) = P(X=x_i, Y) = \begin{cases} \sum_y f_{XY}(x_i, y) & \text{if } X \text{ discrete} \\ \int_Y f_{XY}(x_i, y) dy & \text{if } Y \text{ continuous} \end{cases}$

CONDITIONAL PROBABILITY: $P(X|Y) = \frac{P(X, Y)}{P(Y)}$

BAYES RULE: $P(X|Y) = \frac{P(Y|X) \cdot P(X)}{P(Y)}$

Prior: $P(X)$, encapsulates subjective prior knowledge of the latent variable X

Likelihood: $P(Y|X)$ describes how X and Y are related (= distribution only in y)

Posterior: $P(X|Y)$, expresses what is known about X having observed Y

Evidence: $P(Y)$, marginal likelihood, ensures the posterior is normalized

EXPECTED VALUE: let X be continuous and $g: X \rightarrow \mathbb{R}$ Then $E[g(x)] = \int_X g(x) f_X(x) dx$ (\sum_x if X is discrete), for $X = (X_1, \dots, X_n)$ $E_X[g(x)] = \left(\frac{E_{X_1}[g(x_1)]}{E_{X_2}[g(x_2)]} \right)$

COVARIANCE: let X, Y be random variables Then $\text{cov}[X, Y] = E_{XY}[(X - E_X[X])(Y - E_Y[Y])] = E[XY] - E[X]E[Y]$

N.B.: for $X = (X_1, \dots, X_n), Y = (Y_1, \dots, Y_n)$ $\text{cov}[X, Y] = E[XY^T] - E[X]E[Y^T] = \text{cov}[X, Y]^T$

VARIANCE: let X be a random variable, its variance is defined as $\text{var}[X] := \text{cov}[XX^T] (= \sigma_X^2)$

N.B.: $E[X+Y] = E[X] + E[Y], \text{var}[X+Y] = \text{var}[X] + \text{var}[Y] + \text{cov}[X, Y] + \text{cov}[Y, X]$

CORRELATION: $\text{corr}[X, Y] = \frac{\text{cov}[X, Y]}{\sqrt{\text{var}[X]\text{var}[Y]}} \in [-1, 1]$

EMPIRICAL MEAN: let $X: \Omega \rightarrow \mathbb{R}^D$ and x_1, \dots, x_N be the results of N samplings from X , then the empirical mean of X is given by $\bar{x} := \frac{1}{N} \sum_{i=1}^N x_i \in \mathbb{R}^D$

EMPIRICAL VARIANCE: let $X: \Omega \rightarrow \mathbb{R}^D$ and x_1, \dots, x_N be the results of N samplings from X , then the empirical covariance of X is given by $\Sigma := \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$

N.B.: empirical standard deviation is defined as $\sqrt{\Sigma}$

INDEPENDENCE: $X \perp\!\!\!\perp Y \Leftrightarrow P(X, Y) = P(X)P(Y) \Leftrightarrow P(Y|X) = P(Y) \Leftrightarrow P(X|Y) = P(X)$

N.B.: $X \perp\!\!\!\perp Y \Rightarrow \text{var}[X+Y] = \text{var}[X] + \text{var}[Y] (\Leftrightarrow \text{cov}[X, Y] = 0)$

CONDITIONAL INDEPENDENCE: $X \perp\!\!\!\perp Y|Z \Leftrightarrow P(X, Y|Z) = P(X|Z)P(Y|Z) \forall z \in Z$

MULTIVARIATE GAUSSIAN: let $X \in \mathbb{R}^D$ be a random vector, μ be a mean vector and Σ be a covariance matrix, if $P(X|\mu, \Sigma) = (\frac{1}{2\pi})^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}} \exp(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu))$ Then $X \sim N(\mu, \Sigma)$

N.B.: let $X, Y \sim N(\mu, \Sigma)$ Then $P(X+Y) \sim N(\mu_X + \mu_Y, \Sigma_X + \Sigma_Y)$

PROBABILISTIC MODELS: data = samples from $(X, Y) \sim p(X, Y | \theta)$ where θ = parameters, learn parameters that best fit the data

MAXIMUM LIKELIHOOD ESTIMATION (MLE): probabilistic parameter estimation method, define $f(\theta) := P(y | X\theta) = \prod_{n=1}^N P(y_n | x_n, \theta)$ (= likelihood function), $\theta^* = \arg \max_{\theta} f(\theta)$

N.B.: 1) result is generally computed as $\theta^* = \arg \min_{\theta} -\log(f(\theta)) = L(\theta)$

$$2) L(\theta) = -\log \prod_{n=1}^N P(y_n | x_n, \theta) = -\sum_{n=1}^N \log P(y_n | x_n, \theta)$$

3) MLE is prone to overfitting, parameter values tend to grow large

LINEAR REGRESSION: supervised task, learn a linear function $y = \theta^T x + \epsilon$, where $\epsilon \sim N(0, \sigma^2)$

N.B.: $\epsilon \sim N(0, \sigma^2) \iff P(y | x, \theta) \sim N(\theta^T x, \sigma^2)$, hence parameters can be estimated through MLE where $L(\theta) = \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \theta^T x_i)^2 = \frac{1}{2\sigma^2} (y - X\theta)^T (y - X\theta) = \frac{1}{2\sigma^2} \|y - X\theta\|^2$ (= loss), $\nabla L(\theta) = \frac{1}{\sigma^2} X^T (X\theta - y)$

LEAST SQUARE PROBLEM: minimize $f = \|b - Ax\|^2$ where \exists solutions for $Ax = b$, $\nabla f = 0 \iff A^T A x = A^T b$, furthermore f is convex, thus $\arg \min_x f = \bar{x}$ s.t. $A^T A \bar{x} = A^T b$

N.B.: 1) linear regression can be reduced to LSP ($\arg \min_{\theta} L(\theta) = \arg \min_{\theta} \|y - X\theta\|^2 = \theta^* \text{ s.t. } X^T X \theta^* = X^T y$)

2) if $A \neq \text{full rank}$ the problem can be solved as $x = A^+ b$ (pseudo-inverse)

POLYNOMIAL REGRESSION: it is possible to apply non-linear transformations to the inputs to obtain non-linear outputs: $y = \Phi_k(x)^T \theta + \epsilon = \sum_{k=0}^{K-1} \theta_k \psi_k(x) + \epsilon = \sum_{k=0}^{K-1} \theta_k x^k + \epsilon$ ($K-1$ degree polynomial)

N.B.: 1) let $\Phi_k = \begin{pmatrix} 1 & x_1^k & \dots & x_N^k \end{pmatrix}^T$ be the feature matrix of the model, $L(\theta) = \frac{1}{2\sigma^2} (y - \Phi_k \theta)^T (y - \Phi_k \theta)$, $\theta^* = (\Phi_k^T \Phi_k)^{-1} \Phi_k^T y$, $\nabla L(\theta) = \frac{1}{\sigma^2} \Phi_k^T (\Phi_k \theta - y)$

2) $\Phi_k(x) = \text{Vandermonde}(x, k)^T$

ROOT MEAN SQUARE ERROR (RMSE): measure of error for the model, $\text{RMSE} = \sqrt{\frac{1}{N} \|y - \Phi \theta\|^2} = \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n - \Phi_k^T x_n \theta)^2}$

MAXIMUM A POSTERIORI ESTIMATION (MAP): reduces MLE's overfitting, assume \exists a prior probability $P(\theta)$ on the parameters then $P(\theta | x, y) = \frac{P(y | x, \theta) \cdot P(\theta)}{P(y | x)}$, $\theta^* = \arg \max_{\theta} P(\theta | x, y)$

N.B.: 1) usually $P(\theta) \sim N(0, \sigma^2)$ to limit the parameters' magnitude

2) The evidence $P(y | x)$ is constant w.r.t. θ , negligible

3) assuming $P(\theta) \sim N(0, \sigma^2)$, $L(\theta) = -\log P(\theta | x, y) = -\log P(y | x, \theta) - \log P(\theta) = \frac{1}{2\sigma^2} \|x\theta - y\|^2 + \frac{1}{2\sigma^2} \|\theta\|^2 = \frac{1}{2} \|x\theta - y\|^2 + \frac{1}{2} \|\theta\|^2$ where $\lambda = \frac{\sigma^2}{2}$ (= regularization parameter)

4) $L(\theta)$ convex, $\nabla L(\theta) = 0 \iff (x^T x\theta - x^T y) + \lambda\theta = 0 \iff (x^T x + \lambda I)\theta = x^T y$