

Chapter 2

Machine Learning Fundamentals

Qiang Ji

Table of Contents

- Introduction
- Basic mathematics
 - Probability calculus
 - Linear Algebra
 - Multivariable calculus
 - Gradient-based learning
- Linear Regression
- Linear classification

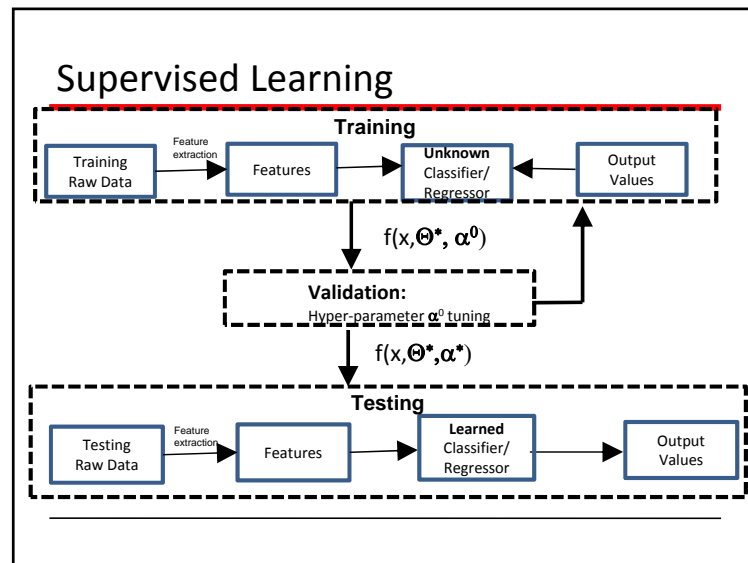
Machine Learning

$$\mathbf{x} \longrightarrow f(\mathbf{x}, \Theta, \alpha) \longrightarrow \mathbf{y}$$

- \mathbf{x} -input, \mathbf{y} -output
- $f()$ the mapping function that maps input to output
- Θ and α are the parameters and hyper-parameters of the mapping function

Types of Learning

- Supervised learning
 - Training data includes both inputs and the corresponding outputs
- Unsupervised learning
 - Training contains only inputs and does not include desired output
- Reinforcement learning
 - Inputs include data and rewards
 - Outputs include actions to take



Basic mathematics

- Probability calculus
- Linear algebra
- Multivariable calculus
- Gradient based learning

Only provide minimally sufficient materials.
Read related references for details.

Probability Calculus

- **Random variable:** a random variable (RV) is a variable whose value is uncertain, depending on its chance as a result of a random process that maps a RV into a specific value.
- Let capital X represent a RV and lower case $x \in \mathcal{X}$ represent a particular value of X , where \mathcal{X} defines the value space
- A RV can be discrete or continuous. A discrete RV assumes a finite set of values while a continuous RV assumes a real value.

Probability Calculus

- The chance of a RV assuming a particular value is measured by its probability, i.e., $P(X=x)$ or $p(x)$ in short, and $0 \leq p(x) \leq 1$
- $P(X)$ represents the probability distribution (pdf) of X and

– For discrete RV,

$$\sum_{x \in \mathcal{X}} p(x) = 1$$

– For continuous RV,

$$\int_{x \in \mathcal{X}} p(x) dx = 1$$

Probability Calculus (cont'd)

- A **random vector** consists of a vector of RVs. We use bold variable to represent a random vector, i.e., $\mathbf{X}=(X_1, X_2, \dots, X_N)^t$, and we use a lower case bold \mathbf{x} to represent a value vector of \mathbf{X} , i.e., $\mathbf{X}=\mathbf{x}$. Both \mathbf{X} and \mathbf{x} are column vectors.
- $p(\mathbf{X})$ represents the probability of the random vector, i.e., the joint probability of all RVs in \mathbf{X} .

Expectation (Mean)

- For a discrete RV X

$$E(X) = \sum_{x \in X} x \cdot p(x)$$
- For a continuous RV X

$$E(X) = \int_x x \cdot p(x) dx$$
- Conditional Expectation

$$E(X | y) = \int_x x \cdot p(x | y) dx$$
- For a random vector $\mathbf{X}=(X_1, X_2, \dots, X_N)^t$

$$E(\mathbf{X})=(E(X_1), E(X_2), \dots, E(X_N))^t$$

Variance

- The variance of a RV X

$$Var(X) = \int [X - E(X)]^2 p(x) dx$$

$$= E[(X - E(X))^2] = E(X^2) - E^2(X)$$
- Standard deviation

$$\sigma_X = \sqrt{Var(X)}$$
- $Var(X|y)=E[(x-E(X|y))^2]=E(X^2|y)-E^2(X|y)$

Covariance and Covariance Matrix

- Covariance of RVs X and Y

$$\sigma_{XY}^2 = E[(X - E(X))(Y - E(Y))]$$

$$= E(XY) - E(X)E(Y)$$
- Covariance Matrix-variance of a random vector $\mathbf{X}=(X_1, X_2, \dots, X_N)^t$

$$\Sigma_X^{N \times N} = E[(\mathbf{X} - E(\mathbf{X}))(\mathbf{X} - E(\mathbf{X}))']$$

Diagonal-variance

Off-diagonal-covariance

Probability Rules

- Product rule
 - Union rule
 - Chain rule
 - Sum rule
 - Conditional probability rule
 - Bayes' rule
-

Product Rule

Given two RVs X and Y , let $p(X,Y)$ represent their joint probability and $p(X|Y)$ represent the conditional probability of X given Y ,

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

$$p(X | Y) = \frac{p(X, Y)}{p(Y)}$$

Union Rule

Given two RVs X and Y , let $p(X+Y)$ represent the probability of their union, i.e., the probability of X or Y .

$$P(X+Y)=p(X)+p(Y)-p(X,Y)$$

if X and Y are mutually exclusive, $p(X,Y)=0$

$$P(X+Y)=p(X)+p(Y)$$

Chain Rule

Given three RVs A, B, C ,

Chain rule :

$$p(A, B, C) = p(A) p(B | A) p(C | A, B)$$

Conditional chain rule:

$$p(A, B, C | D, E) = p(A | D, E) p(B | A, D, E) p(C | A, B, D, E)$$

Any assumptions? Chain rule for N variables X_1, X_2, \dots, X_N ?

Sum Rule

- Sum rule via marginalization

$$P(X) = \sum_Y P(X, Y) = \sum_Y \sum_Z P(X, Y, Z)$$

- Sum rule via conditional probability

$$P(X) = \sum_Y P(X | Y) p(Y)$$

$$P(X | Y) = \sum_Z P(X | Y, Z) p(Z | Y)$$

Bayes' Rule

$$p(X | Y) = \frac{\overset{\text{prior probability of X}}{p(X)} \overset{\text{Likelihood of X}}{p(Y | X)}}{p(Y)}$$

$p(Y)$ is a normalization constant $p(Y) = \sum_X p(X) p(Y | X)$

Independence

- If X and Y are marginally independent, then

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

$$p(X | Y) = ?$$

$$E(XY) = ?$$

$$\text{Cov}(X, Y) = ?$$

We denote it as $X \perp Y$

Conditional Independence

- For three RVs, X , Y , and Z , if X and Y are independent, given Z , we have

$$P(X | Y, Z) = P(X | Z)$$

$$p(X, Y | Z) = ?$$

$$E(X, Y | Z) = ?$$

We denote it as

$$X \perp Y | Z$$

Probability Distributions

- Probability distribution can be **discrete** or **continuous**.
- Probability distribution for one RV, i.e., **univariate** distribution or probability distribution for multiple RVs, i.e., **multivariate** or joint distribution

Discrete Probability Distributions

- Uniform distribution $X \in \{1, 2, \dots, K\}$
 $X \sim \text{Uniform}(x | K) \quad P(X = k | K) = \frac{1}{K}$
- Bernoulli distribution for a binary RV $X \in \{0, 1\}$
 $X \sim \text{Ber}(x | \theta), P(X = 1) = \theta, p(X = 0) = 1 - \theta$
- Binomial distribution
 ➤ Let Y be an integer RV that represents the number of times for $X=1$ for N Bernoulli trials, with $p(X = 1) = \theta$
 $Y \sim \text{Bin}(y | N, \theta) \quad P(Y = n_1 | N, \theta) = \binom{N}{n_1} \theta^{n_1} (1 - \theta)^{N - n_1}$

Discrete Probability Distributions

- Multinomial distribution for K discrete RVs, Y_1, Y_2, \dots, Y_K , where Y_k is an integer that represents the number of times a discrete RV $X \in \{1, 2, \dots, K\}$ equal to k out of a total of N trials of X . Let $\mathbf{Y} = (Y_1, Y_2, \dots, Y_K)$

$$Y \sim \text{Mult}(y_1, y_2, \dots, y_K | N, \theta_1, \theta_2, \dots, \theta_K) = \frac{N!}{y_1! y_2! \dots y_K!} \theta_1^{y_1} \theta_2^{y_2} \dots \theta_K^{y_K}$$

where $p(X = k) = \theta_k$

Binomial distribution is a special case of multinomial distribution, where $K=2$

Continuous Probability Distributions

- Gaussian distribution
 - Unary: $p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
 - Multivariate $\mathbf{X} \sim N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}} e^{-\frac{(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}{2}}$
- Dirichlet distribution
 - Continuous multivariate probability distributions for K RVs X_1, X_2, \dots, X_K , where $X_i \in (0, 1)$ and $\sum_{i=1}^K X_i = 1$
 $\mathbf{X} \sim \text{Dir}(x_1, x_2, \dots, x_K | \alpha_1, \alpha_2, \dots, \alpha_K) = p(x_1, x_2, \dots, x_K | \alpha_1, \alpha_2, \dots, \alpha_K) = \frac{1}{B(\alpha_1, \alpha_2, \dots, \alpha_K)} \prod_{i=1}^K x_i^{\alpha_i - 1}$
 where $B(\alpha_1, \alpha_2, \dots, \alpha_K) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^K \alpha_i)}$ and $\Gamma(n) = (n-1)!$ gamma function

Linear Algebra

- A scalar is a single number. It can be real (\mathbb{R}), integer (\mathbb{N}), boolean, binary, etc. . It is represent by a lower case symbol: x, n, i, j

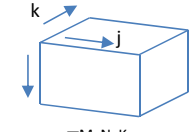
$x \in \mathbb{R}$ – a real scalar $x \in \mathbb{N}$ – an integer scalar

- A vector is a 1-D array of scalars, represented by a bold upper case letter in column vector

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{pmatrix} \quad \mathbf{x} \in \mathbb{R}^N \quad \text{a vector of real numbers}$$

Linear Algebra

- A matrix \mathbf{A} (i,j) is a 2D array of scalars, represented by a uppercase bold letter

$$\mathbf{A}^{M \times N} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1N} \\ a_{21} & a_{22} & \dots & a_{2N} \\ \dots & \dots & \dots & \dots \\ a_{M1} & a_{M2} & \dots & a_{MN} \end{pmatrix} \quad \mathbf{A} \in \mathbb{R}^{M \times N} \quad i \quad j$$


$\mathbf{T}^{M \times N \times K}$

- A tensor \mathbf{T} is a N-D dimensional array of scalars $\mathbf{T}(i,j,k,..)$. A tensor becomes a scalar, vector, and matrix when $N=0, 1$, and 2 . When $N=3$, tensor \mathbf{T} is a 3D matrix $\mathbf{T}(i,j,k)$.

Vector Norms

$$\mathbf{X} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{pmatrix} \quad \|\mathbf{X}\|_p = \left(\sum_{i=1}^N |x_i|^p \right)^{\frac{1}{p}}$$

$$L_1 : \|\mathbf{X}\|_1 = \sum_{i=1}^N |x_i|$$

$$L_2 : \|\mathbf{X}\|_2 = \left(\sum_{i=1}^N x_i^2 \right)^{\frac{1}{2}} = \sqrt{x_1^2 + x_2^2 + \dots + x_N^2}, \text{ the Euclidean distance}$$

$$L_\infty : \|\mathbf{X}\|_\infty = \max |x_i|$$

$$L_0 = \#(i | x_i \neq 0) - \text{the cardinality of } \mathbf{X}$$

Matrix Norms

$$\mathbf{A}^{M \times N} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1N} \\ a_{21} & a_{22} & \dots & a_{2N} \\ \dots & \dots & \dots & \dots \\ a_{M1} & a_{M2} & \dots & a_{MN} \end{pmatrix}$$

$$L_1 : \|\mathbf{A}\|_1 = \max_{1 \leq j \leq M} \sum_{i=1}^N |a_{ij}| = \max \left(\sum_{i=1}^N |a_{i1}|, \sum_{i=1}^N |a_{i2}|, \dots, \sum_{i=1}^N |a_{iM}| \right)$$

$$L_2 : \|\mathbf{A}\|_2 = \left(\sum_{i=1}^M \sum_{j=1}^N a_{ij}^2 \right)^{\frac{1}{2}}$$

Trace and Determinant

- Trace of matrix
 - Sum of the diagonal elements, i.e.,

$$Tr(\mathbf{A}) = \sum_i a_{i,i}$$

- Determinant of A
 - Product of all eigenvalues of A, i.e.,
- Rank of A
 - The number of independent rows or columns

Eigenvalues and Eigenvectors

- An eigenvector of a square matrix $\mathbf{A}^{N \times N}$ is a non-zero vector \mathbf{v} such that multiplication of \mathbf{A} by \mathbf{v} alters only the scale of \mathbf{v} , i.e.,

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

where λ is called the eigenvalue of A.

- A has N eigenvectors (\mathbf{v}_n) and N eigenvalues (λ_n)

Eigen-decomposition

- Given its eigenvalues λ_n and eigenvectors \mathbf{v}_n of a matrix, a matrix \mathbf{A} can be decomposed into

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1} \quad \mathbf{V}^{N \times N} = [\mathbf{v}_1 \mathbf{v}_2 \dots \mathbf{v}_N] \quad \mathbf{\Lambda}^{N \times N} = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \dots \\ & & & \lambda_N \end{bmatrix}$$

Note \mathbf{V} is orthonormal, i.e., $\mathbf{V}^{-1} = \mathbf{V}^t$ if A is real and symmetric.

Singular Value Decomposition (SVD)

- A matrix $\mathbf{A}^{M \times N}$ can be decomposed into the products of three matrices

$$\mathbf{A}^{M \times N} = \mathbf{U}^{M \times M} \mathbf{D}^{M \times N} \mathbf{V}^{N \times N}$$

where U and V are orthonormal matrices ($\mathbf{U}^{-1} = \mathbf{U}^t$) and D is a diagonal matrix, whose values are called singular values

SVD v.s. Eigen-decomposition

- For a square symmetric matrix $\mathbf{A}^{M \times M}$, its SVD is

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}^t$$

- Its eigen decomposition is

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^{-1}$$

- Hence, $\mathbf{U} = \mathbf{V}$ and $\mathbf{D} = \mathbf{\Lambda}$

SVD v.s. Eigen-decomposition

- For a non-square matrix $\mathbf{A}^{M \times N}$, its SVD is

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^t$$

- The columns of \mathbf{V} are eigenvectors of $\mathbf{A}^t \mathbf{A}$.
- The columns of \mathbf{U} are eigenvectors $\mathbf{A} \mathbf{A}^t$
- The non-zero singular values of \mathbf{D} are the square roots of the nonzero eigenvalues of $\mathbf{A} \mathbf{A}^t$ (or $\mathbf{A}^t \mathbf{A}$)

SVD Applications

- The rank of \mathbf{A} is the number of non-zero singular values.

- Computing matrix inverse

$$\mathbf{A}^{-1} = (\mathbf{U} \mathbf{D} \mathbf{V}^t)^{-1} = \mathbf{V} \mathbf{D}^{-1} \mathbf{U}^t$$

- Computing eigen vectors

- For a symmetric matrix \mathbf{A} , its SVD decomposition

$$\mathbf{A} = \mathbf{V} \mathbf{D} \mathbf{V}^t$$

- The columns of \mathbf{V} are eigenvectors of \mathbf{A} and the singular values in \mathbf{D} are the eigenvalues of \mathbf{A}

System of Linear Equations

- Let $\mathbf{x}^{N \times 1}$ be a unknown vector, $\mathbf{A}^{M \times N}$ and $\mathbf{b}^{M \times 1}$ are given. Find \mathbf{x} by minimizing

$$(\mathbf{A} \mathbf{x} - \mathbf{b})^t (\mathbf{A} \mathbf{x} - \mathbf{b})$$

- If $M \geq N$ and the rows are independent, \mathbf{x} has a unique solution

$$\mathbf{x} = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t \mathbf{b}$$

where $(\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t$ is called pseudo-inverse of \mathbf{A} can be computed via SVD

- If $M < N$, \mathbf{x} has multiple solutions

Principal Component Analysis (PCA)

- Given data $\mathbf{D}=[\mathbf{x}[m]], m=1,2,\dots,M$, compute the covariance matrix Σ for samples in \mathbf{D}
- Perform SVD on $\Sigma=UDV^t$, and obtain the eigen vectors of from columns of V matrix
- Order the eigen vectors in a descending order according to their eigen values
- Select top K eigen vectors and form a project matrix M , whose rows are the top K selected eigen vectors
- Multiply each input \mathbf{x} by M to produce the projected data \mathbf{y} , i.e., $\mathbf{y}=\mathbf{M}\mathbf{x}$

Multivariable (Vector) Calculus

- Involves differentiation and integration of vectors and matrices
- Derivatives with vector
 - Scalar by vector
 - Vector by scalar
 - Vector by vector
- Derivatives with Matrices
 - Scalar by matrix
 - Vector by matrix
- Follow denominator layout convention

Derivatives with Vectors

Let $\mathbf{X}^{N \times 1}$ and $\mathbf{Y}^{M \times 1}$ be column vectors and z be a scalar function of $z(\mathbf{X})$, following the denominator layout

- Scalar by vector

$$\frac{\partial z}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial z}{\partial X_1} \\ \frac{\partial z}{\partial X_2} \\ \dots \\ \frac{\partial z}{\partial X_N} \end{pmatrix}$$
- Vector by scalar

$$\frac{\partial \mathbf{X}^{1 \times N}}{\partial z} = \begin{pmatrix} \frac{\partial X_1}{\partial z} & \frac{\partial X_2}{\partial z} & \dots & \frac{\partial X_N}{\partial z} \end{pmatrix}$$

Vector by Vector

- Vector by vector

$$\left(\frac{\partial \mathbf{X}^{N \times 1}}{\partial \mathbf{Y}^{M \times 1}} \right)^{M \times N} = \begin{pmatrix} \frac{\partial X_1}{\partial Y_1} & \frac{\partial X_2}{\partial Y_1} & \dots & \frac{\partial X_N}{\partial Y_1} \\ \frac{\partial X_1}{\partial Y_2} & \frac{\partial X_2}{\partial Y_2} & \dots & \frac{\partial X_N}{\partial Y_2} \\ \dots & \dots & \dots & \dots \\ \frac{\partial X_1}{\partial Y_M} & \frac{\partial X_2}{\partial Y_M} & \dots & \frac{\partial X_N}{\partial Y_M} \end{pmatrix}$$
- Let $\mathbf{V}^{M \times 1}$ be column vector function $\mathbf{X}^{N \times 1}$ and $A^{K \times M}$ be a matrix that is NOT a function of \mathbf{X}

$$\frac{\partial (A\mathbf{V})^{N \times K}}{\partial \mathbf{X}} = \frac{\partial \mathbf{V}}{\partial \mathbf{X}} A^t$$

Scalar Function by Vector

Let $\mathbf{U}^{M \times 1}(\mathbf{X})$ and $\mathbf{V}^{M \times 1}(\mathbf{X})$ be column vectors, and both are a function $\mathbf{X}^{N \times 1}$.

$$\frac{\partial(\mathbf{U}'\mathbf{V})}{\partial \mathbf{X}} = \frac{\partial \mathbf{U}}{\partial \mathbf{X}} \mathbf{V} + \frac{\partial \mathbf{V}}{\partial \mathbf{X}} \mathbf{U}$$

Scalar Function by Vector

Let $\mathbf{U}^{K \times 1}$ and $\mathbf{V}^{M \times 1}$ be column vectors, and both are a function $\mathbf{X}^{N \times 1}$, $\mathbf{A}^{K \times M}$ be a matrix that is NOT a function of \mathbf{X}

$$\frac{\partial(\mathbf{U}'\mathbf{A}\mathbf{V})}{\partial \mathbf{X}} = \frac{\partial \mathbf{V}}{\partial \mathbf{X}} \mathbf{A}' \mathbf{U} + \frac{\partial \mathbf{U}}{\partial \mathbf{X}} \mathbf{A} \mathbf{V}$$

Scalar Function by Vector

Given $\mathbf{A}^{M \times N}$, $\mathbf{b}^{M \times 1}$, $\mathbf{C}^{M \times K}$, $\mathbf{D}^{K \times N}$, $\mathbf{e}^{K \times 1}$ that are Not a function $\mathbf{X}^{N \times 1}$,

$$\begin{aligned} \frac{\partial(\mathbf{A}\mathbf{X} + \mathbf{b})' \mathbf{C}(\mathbf{D}\mathbf{X} + \mathbf{e})}{\partial \mathbf{X}} &= \frac{\partial(\mathbf{D}\mathbf{X} + \mathbf{e})}{\partial \mathbf{X}} \mathbf{C}'(\mathbf{A}\mathbf{X} + \mathbf{b}) + \frac{\partial(\mathbf{A}\mathbf{X} + \mathbf{b})}{\partial \mathbf{X}} \mathbf{C}(\mathbf{D}\mathbf{X} + \mathbf{e}) \\ &= \mathbf{D}' \mathbf{C}'(\mathbf{A}\mathbf{X} + \mathbf{b}) + \mathbf{A}' \mathbf{C}(\mathbf{D}\mathbf{X} + \mathbf{e}) \end{aligned}$$

Derivatives with Matrices

Let $\mathbf{A}^{M \times N}$ be a matrix and z be a scalar function of \mathbf{A}

- Scalar by matrix

$$\left(\frac{\partial z}{\partial \mathbf{A}} \right)^{1 \times MN} = \left(\frac{\partial z}{\partial \mathbf{A}[1]1} \quad \frac{\partial z}{\partial \mathbf{A}[1]2} \quad \dots \quad \frac{\partial z}{\partial \mathbf{A}[1]N} \right) = \left(\frac{\partial z}{\partial \mathbf{A}[2]1} \quad \frac{\partial z}{\partial \mathbf{A}[2]2} \quad \dots \quad \frac{\partial z}{\partial \mathbf{A}[2]N} \right) = \dots = \left(\frac{\partial z}{\partial \mathbf{A}[M]1} \quad \frac{\partial z}{\partial \mathbf{A}[M]2} \quad \dots \quad \frac{\partial z}{\partial \mathbf{A}[M]N} \right)$$

where $\mathbf{A}[i]j$ - i -th column of \mathbf{A} , $\mathbf{A}[i]j$ - j -th row of \mathbf{A}

- Let $\mathbf{u}^{M \times 1}$ and $\mathbf{v}^{N \times 1}$ are not function of \mathbf{A}

$$\frac{\partial \mathbf{u}' \mathbf{A} \mathbf{v}}{\partial \mathbf{A}} = \mathbf{u} \mathbf{v}'$$

Derivatives with Matrices

- Vector by matrix-Tensor . Let $\mathbf{X}^{K \times 1}$ be a vector

$$\left(\frac{\partial \mathbf{X}^{K \times 1}}{\partial \mathbf{A}^{M \times N}} \right)^{M \times N \times K} = \left(\frac{\partial \mathbf{X}}{\partial \mathbf{A}[\cdot][1]} \frac{\partial \mathbf{X}}{\partial \mathbf{A}[\cdot][2]} \cdots \frac{\partial \mathbf{X}}{\partial \mathbf{A}[\cdot][N]} \right) = \left(\frac{\partial \mathbf{X}[1]}{\partial \mathbf{A}} \frac{\partial \mathbf{X}[2]}{\partial \mathbf{A}} \cdots \frac{\partial \mathbf{X}[K]}{\partial \mathbf{A}} \right)$$

$$= \begin{pmatrix} \frac{\partial \mathbf{X}}{\partial \mathbf{A}[1][\cdot]} \\ \frac{\partial \mathbf{X}}{\partial \mathbf{A}[2][\cdot]} \\ \cdots \\ \frac{\partial \mathbf{X}}{\partial \mathbf{A}[M][\cdot]} \end{pmatrix} = \begin{pmatrix} \frac{\partial \mathbf{X}[1]}{\partial \mathbf{A}[1][\cdot]} & \frac{\partial \mathbf{X}[2]}{\partial \mathbf{A}[1][\cdot]} & \cdots & \frac{\partial \mathbf{X}[K]}{\partial \mathbf{A}[1][\cdot]} \\ \frac{\partial \mathbf{X}[1]}{\partial \mathbf{A}[2][\cdot]} & \frac{\partial \mathbf{X}[2]}{\partial \mathbf{A}[2][\cdot]} & \cdots & \frac{\partial \mathbf{X}[K]}{\partial \mathbf{A}[2][\cdot]} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial \mathbf{X}[1]}{\partial \mathbf{A}[M][\cdot]} & \frac{\partial \mathbf{X}[2]}{\partial \mathbf{A}[M][\cdot]} & \cdots & \frac{\partial \mathbf{X}[K]}{\partial \mathbf{A}[M][\cdot]} \end{pmatrix}$$

Derivatives with Matrices

- Tensor vector multiplication

Let \mathbf{Y} be a $K \times 1$ vector

$$\frac{\partial \mathbf{X}^{M \times N \times K}}{\partial \mathbf{A}} \mathbf{Y}^{K \times 1} = \left(\sum_{k=1}^K \frac{\partial \mathbf{X}_k}{\partial \mathbf{A}} \mathbf{Y}_k \right)^{M \times N}$$

- Tensor matrix multiplication

Let \mathbf{B} be a $K \times D$ matrix

$$\frac{\partial \mathbf{X}^{M \times N \times K}}{\partial \mathbf{A}} \mathbf{B}^{K \times D} = \sum_{k=1}^K \frac{\partial \mathbf{X}_k}{\partial \mathbf{A}} \mathbf{B}[k][\cdot]$$

Derivatives with Matrices

- Derivative of Determinant of A

$$\frac{\partial |\mathbf{A}|}{\partial \mathbf{A}} = |\mathbf{A}| \mathbf{A}^{-T}$$

Parameter Learning

$$\mathbf{x} \longrightarrow f(\mathbf{x}, \Theta) \longrightarrow \mathbf{y}$$

- Parameter learning is to learn the mapping function parameters Θ , given M i.i.d training data $\mathbf{D} = \{\mathbf{X}[m], \mathbf{Y}[m]\}$, where $m=1, 2, \dots, M$
- Parameter learning is often formulated as finding Θ to minimize certain loss function $L(\mathbf{D}; \Theta)$.

Parameter Learning (cont'd)

- Let $l(X[m], y[m])$ be the loss function for each sample, $L(\mathbf{D}; \Theta)$ is defined below
 - Unregularized loss function (average loss)
 - $L(\mathbf{D}; \Theta) = \frac{1}{M} \sum_{m=1}^M l(X[m], y[m], \Theta)$,
 - Regularized loss function
 - $L(\mathbf{D}; \Theta) = \frac{1}{M} \sum_{m=1}^M l(X[m], y[m], \Theta) + \lambda R(\Theta)$, $R(\Theta)$ is a regularization term
- Given the loss function, Θ is solved by minimizing the loss function, i.e.,

$$\Theta^* = \arg \min_{\Theta} L(\mathbf{D}; \Theta)$$

Gradient-based Parameter Learning

For some loss function, analytic closed form solution exists. Parameters Θ can be solved by computing the gradient of the loss function with respect to Θ and setting it to zero

$$\nabla_{\Theta} L(\mathbf{D}; \Theta) = \frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta} = 0$$

Gradient based Parameter Learning

- For many problems, analytic solution may not exist, **gradient descent** is often used to iteratively estimate Θ

$$\Theta^t = \Theta^{t-1} - \eta_t \left. \frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta} \right|_{\Theta = \Theta^{t-1}}$$

where Θ^t is the estimate at t th iteration, η is the learning rate that needs be manually tuned. Θ needs be initialized to initial value Θ^0

- Gradient descent method is the dominant learning technique for deep learning

Gradient Descent

- Descent the mountain iteratively. At each iteration, the best direction to descend is the negative of the gradient
- For convex loss function, descent will converge to global minimum independent of initialization.
- For non-convex loss function, descent may stuck in a local minimum
- In practice, local minimum does not seem to be a problem-a theoretical mystery!



Gradient based Parameter Learning

- For non-regularized loss function, gradient of the parameters can be computed as follows

$$\frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta} = \frac{1}{M} \sum_{m=1}^M \frac{\partial l(\mathbf{x}[m], y[m], \Theta)}{\partial \Theta}$$

- For regularized loss function,

$$\frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta} = \frac{1}{M} \sum_{m=1}^M \frac{\partial l(\mathbf{x}[m], y[m], \Theta)}{\partial \Theta} + \lambda \frac{\partial R(\Theta)}{\partial \Theta}$$

Stochastic Gradient Descent (SGD)

- When the number of samples (M) is large, computing the gradient of Θ using all samples can be very slow.

$$\frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta} = \frac{1}{M} \sum_{m=1}^M \frac{\partial l(\mathbf{x}[m], y[m], \Theta)}{\partial \Theta}$$

- Stochastic Gradient Descent (SGD) can approximately compute the gradient

Stochastic Gradient Descent

Mini-batch approximation

- Divide \mathbf{D} into K batches- D_1, D_2, \dots, D_K with the batch size S ranging from 2 to 100.
- At each iteration t, randomly select a batch D_k and compute the gradient

$$\frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta} = \frac{1}{S} \sum_{\mathbf{x}[m], y[m] \in D_k} \frac{\partial l(\mathbf{x}[m], y[m], \Theta)}{\partial \Theta}$$

$$\Theta' = \Theta^{t-1} - \eta_t \left. \frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta} \right|_{\Theta = \Theta^{t-1}}$$

- The size of the batch (a hyper-parameter) is independent of the number of samples

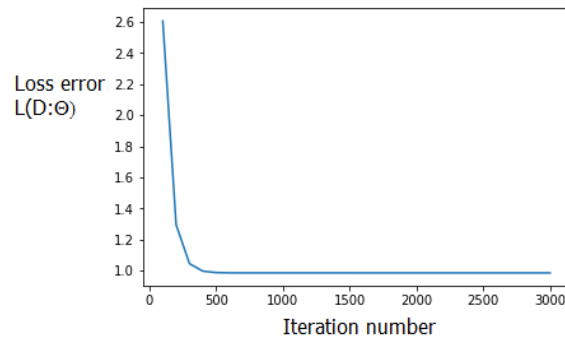
Convergence of GD/SGD

- The gradient descent method iterates from an initial point until convergence
- The convergence can be measured by the change of
 - the estimated parameters Θ or
 - the magnitude of the gradient $|\nabla L(\mathbf{D}; \Theta)|$ or
 - the loss function value

When the change is below a threshold, the iteration can stop.

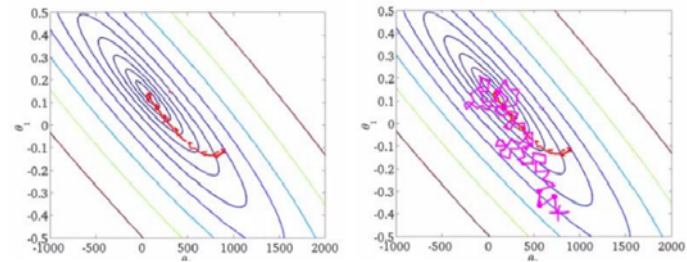
- The iteration can also be stopped when the maximum iteration number is reached.

Convergence of GD/SGD



GD versus SGD

- Batch gradient descent vs. Stochastic Gradient Descent



Source: <https://www.cs.cmu.edu/~yuxiangw/docs/SSGD.pdf>

Regression

The goal of regression is to predict the value of one or more continuous target output y (the *regress and or dependent variable*) values given the value of an input feature vector \mathbf{x}

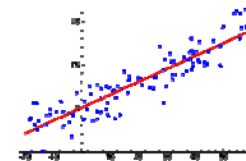
$$f: \mathbf{x} \in \mathbb{R}^N \rightarrow y \in \mathbb{R}^1$$

Linear Regression

- Linear regression means the output y is a linear function of the parameters \mathbf{w} of the regression model, i.e.,

$$y = \mathbf{w}^T \mathbf{x} + w_0$$

where \mathbf{w} is a vector of weight coefficients and w_0 a bias



Linear Regression (cont'd)

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}'\mathbf{x} + w_0 = [\mathbf{x}' \ 1] \begin{bmatrix} \mathbf{w} \\ w_0 \end{bmatrix} = \mathbf{X}'\boldsymbol{\Theta}$$

where $\mathbf{x}^{N \times 1} = (x_1, x_2, \dots, x_N)^t$ is input vector, y (a real scalar), the output variable, and $\mathbf{w} = \{w_1, w_2, \dots, w_N\}$ is the model parameters, and w_0 is the bias.

Given $\boldsymbol{\Theta} = \begin{bmatrix} \mathbf{w} \\ w_0 \end{bmatrix}$, linear regression is to predict the value of y for a new input \mathbf{x}

Linear Regression Learning

Given paired training data $\mathbf{D} = \{\mathbf{x}[m], y[m]\}$, $m=1, 2, \dots, M$, learn the parameters $\boldsymbol{\Theta}$ of the regression function $f(\mathbf{x}, \boldsymbol{\Theta})$ by minimizing some loss functions, i.e.,

$$\boldsymbol{\Theta}^* = \arg \min_{\boldsymbol{\Theta}} L(\mathbf{D} : \boldsymbol{\Theta})$$

where $L()$ is a loss function. Types of loss functions:

1) Mean squared errors (MSE)

$$l(\mathbf{x}[m], y[m]) = (\mathbf{x}^t[m]\mathbf{w} + w_0 - y[m])^2$$

2) Negative Log Conditional Likelihood (-LCL) (also called cross-entropy)

$$l(\mathbf{x}[m], y[m]) = -\log p(y[m]|\mathbf{x}[m], \boldsymbol{\Theta})$$

Mean Squared Errors

Let $\boldsymbol{\Theta} = \begin{bmatrix} \mathbf{w} \\ w_0 \end{bmatrix}$ be the model parameters

and $\mathbf{x} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$ be the homogeneous representation of the input

$$\begin{aligned} L_{MSE}(\boldsymbol{\Theta} : \mathbf{D}) &= \frac{1}{M} \sum_{m=1}^M (\mathbf{w}'\mathbf{x}[m] + w_0 - y[m])^2 \\ &= \frac{1}{M} \sum_{m=1}^M (\mathbf{X}'[m]\boldsymbol{\Theta} - y[m])^2 \\ &= \frac{1}{M} \begin{bmatrix} \mathbf{x}^t[1] \\ \mathbf{x}^t[2] \\ \vdots \\ \mathbf{x}^t[M] \end{bmatrix} \boldsymbol{\Theta} - \begin{bmatrix} y[1] \\ y[2] \\ \vdots \\ y[M] \end{bmatrix} \\ &\quad \begin{matrix} \mathbf{A} & \mathbf{Y} \end{matrix} \\ &= \frac{1}{M} (\mathbf{A}\boldsymbol{\Theta} - \mathbf{Y})'(\mathbf{A}\boldsymbol{\Theta} - \mathbf{Y}) \end{aligned}$$

Least-squares Estimation

$$\boldsymbol{\Theta}^* = \arg \min_{\boldsymbol{\Theta}} L_{MSE}(\mathbf{D} : \boldsymbol{\Theta})$$

$$\frac{\partial L_{MSE}(\mathbf{D} : \boldsymbol{\Theta})}{\partial \boldsymbol{\Theta}} = 0$$

$$\boldsymbol{\Theta}^* = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t \mathbf{Y}$$

Negative Log Conditional Likelihood

$$y = \mathbf{w}'\mathbf{x} + w_0 + \zeta$$

$$\zeta \sim N(0, \sigma^2)$$

$$p(y | \mathbf{x}) \sim N(\mathbf{w}'\mathbf{x} + w_0, \sigma^2)$$

$$\begin{aligned} L_{-LCL}(D : \Theta) &= -\sum_{m=1}^M \log p(y[m] | \mathbf{x}[m]) && \text{cross-entropy} \\ &= -\sum_{m=1}^M \log \left\{ \frac{1}{\sqrt{(2\pi)\sigma}} e^{-\frac{(y[m] - \mathbf{w}'\mathbf{x}[m] - w_0)^2}{2\sigma^2}} \right\} \\ &= \sum_{m=1}^M -\frac{(y[m] - \mathbf{w}'\mathbf{x}[m] - w_0)^2}{2\sigma^2} + M \log \sigma + \frac{M}{2} \log(2\pi) \\ &\approx \frac{(\mathbf{A}\Theta - \mathbf{Y})'(\mathbf{A}\Theta - \mathbf{Y})}{2\sigma^2} + M \log \sigma \end{aligned}$$

Negative Log Conditional Likelihood

$$L_{-LCL}(D : \Theta) = \frac{(\mathbf{A}\Theta - \mathbf{Y})'(\mathbf{A}\Theta - \mathbf{Y})}{2\sigma^2} + M \log \sigma$$

- It is clear that the mean squared loss function is a special case of the negative log conditional likelihood function with σ treated as constant
 - σ is an additional parameter that can also be learnt
-

Maximum Likelihood Estimation

$$\Theta^* = \arg \min_{\Theta} L_{-LCL}(D : \Theta)$$

$$= \arg \max_{\Theta} LCL(D : \Theta)$$

$$\frac{\partial L_{-LCL}(D : \Theta)}{\partial \Theta} = 0 \Rightarrow$$

$$\Theta^* = (\mathbf{A}'\mathbf{A})^{-1} \mathbf{A}'\mathbf{Y}$$

$$\sigma^2 = \frac{(\mathbf{A}\Theta^* - \mathbf{Y})'(\mathbf{A}\Theta^* - \mathbf{Y})}{M}$$

Learning with Regularization

- To control over-fitting, a regularization term is added to the loss function, creating a new target function

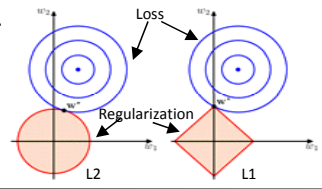
$$L(\Theta : \mathbf{D}) + \lambda R(\Theta)$$

Loss function + Regularization term

- $R(\Theta)$ is typically squared L2 or L1 norm
 - Squared L2 norm $R(\Theta) = \|\Theta\|_2^2 = \Theta'\Theta$
 - L1 norm $R(\Theta) = \|\Theta\|_1$
 - λ is a hyper-parameter that determines the relative weight of the two terms
-

L2 versus L1 norm

- L1 normal (lasso) tends to generate sparser solutions by forcing parameters to become zero. It is convex but non-differentiable at zero.
- L2 norm (quadratic norm) produces small value (close to zero) parameters. It is differentiable and the regularized function remains convex. It is widely used.



Regression Learning with L1 Norm

- With a MSE loss function and a quadratic regularizer $g(\Theta) = \Theta^T \Theta$, the objective function becomes

$$L_{MSEL2}(D : \Theta) = (A\Theta - Y)^T (A\Theta - Y) + \lambda \Theta^T \Theta$$

$$\frac{\partial L_{MSEL2}(D : \Theta)}{\partial \Theta} = 2A^T (A\Theta - Y) + 2\lambda \Theta = 0$$

$$\Theta = (A^T A + \lambda I)^{-1} A^T Y$$

Note I is an identity matrix

Regression Learning with L2 Norm

- With a -LCL loss function and a quadratic regularizer $g(\Theta) = \Theta^T \Theta$, the objective function becomes

$$L_{-LCLL2}(D : \Theta) = \frac{(A\Theta - Y)^T (A\Theta - Y)}{2\sigma^2} + M \log \sigma + \lambda \Theta^T \Theta$$

$$\frac{\partial L_{-LCLL2}(D : \Theta)}{\partial \Theta} = \frac{A^T (A\Theta - Y)}{\sigma^2} + 2\lambda \Theta = 0$$

$$\Theta = (A^T A + 2\lambda \sigma^2 M I)^{-1} A^T Y$$

- $2\lambda \sigma^2$ can be treated as one composite hyper-parameter

Regression Learning with L1 Regularization

- With a MSE loss function and a L1 norm $R(\Theta) = |\Theta|_1$, the objective function becomes

$$\begin{aligned} L_{MSEL1}(D : \Theta) &= (A\Theta - Y)^T (A\Theta - Y) + \lambda |\Theta|_1 \\ &= (A\Theta - Y)^T (A\Theta - Y) + \lambda \sum_{i=1}^N |\Theta_i| \end{aligned}$$

Regression Learning with L1 Regularization

$$\begin{aligned}\frac{\partial \text{MSEL1}(\mathbf{D}; \boldsymbol{\Theta})}{\partial \boldsymbol{\Theta}} &= 2\mathbf{A}'(\mathbf{A}\boldsymbol{\Theta} - \mathbf{Y}) + \lambda \frac{\partial \|\boldsymbol{\Theta}\|_1}{\partial \boldsymbol{\Theta}} \\ &= 2\mathbf{A}'(\mathbf{A}\boldsymbol{\Theta} - \mathbf{Y}) + \lambda \begin{pmatrix} \frac{\partial \|\boldsymbol{\Theta}\|_1}{\partial \Theta_1} \\ \frac{\partial \|\boldsymbol{\Theta}\|_1}{\partial \Theta_2} \\ \dots \\ \frac{\partial \|\boldsymbol{\Theta}\|_1}{\partial \Theta_N} \end{pmatrix} = 2\mathbf{A}'(\mathbf{A}\boldsymbol{\Theta} - \mathbf{Y}) + \lambda \begin{pmatrix} \text{sign}(\Theta_1) \\ \text{sign}(\Theta_2) \\ \dots \\ \text{sign}(\Theta_N) \end{pmatrix}\end{aligned}$$

$$\text{sign}(\Theta_i) = \begin{cases} 1 & \text{if } \Theta_i > 0 \\ -1 & \text{if } \Theta_i < 0 \\ c \in [-1, +1] & \text{if } \Theta_i = 0, c \text{ is sometimes chosen to be zero} \end{cases}$$

Hence, $\frac{\partial \|\boldsymbol{\Theta}\|_1}{\partial \boldsymbol{\Theta}}$ equals 2^N vectors of +1s or -1s, and it is nondifferentiable at $\Theta_i = 0$

Subgradient descent method can be used.