

CS6720 --- Pattern Recognition

Review of Prerequisites in Math and Statistics

Prepared by
Li Yang

Based on
Appendix chapters of
Pattern Recognition, 4th Ed.
by S. Theodoridis and K. Koutroumbas
and figures from
Wikipedia.org

1

Probability and Statistics

- ❖ **Probability** $P(A)$ of an event A : a real number between 0 to 1.
- ❖ **Joint probability** $P(A \cap B)$: probability that both A and B occurs in a single experiment.
 $P(A \cap B) = P(A)P(B)$ if A and B are **independent**.
- ❖ Probability $P(A \cup B)$ of union of A and B: either A or B occurs in a single experiment.
 $P(A \cup B) = P(A) + P(B)$ if A and B are mutually exclusive.

- ❖ **Conditional probability:**

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

- ❖ Therefore, the **Bayes rule:**

$$P(A|B)P(B) = P(B|A)P(A) \text{ and } P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- ❖ **Total probability:** let A_1, \dots, A_m such that $\sum_{i=1}^m P(A_i) = 1$ then

$$P(B) = \sum_{i=1}^m P(B|A_i)P(A_i)$$

2

- ❖ **Probability density function (pdf):** $p(x)$ for a continuous random variable x

$$P(a \leq x \leq b) = \int_a^b p(x) dx$$

Total and conditional probabilities can also be extended to pdf's.

- ❖ **Mean and Variance:** let $p(x)$ be the pdf of a random variable x

$$E[x] = \int_{-\infty}^{\infty} xp(x) dx, \text{ and } \sigma^2 = \int_{-\infty}^{\infty} (x - E[x])^2 p(x) dx$$

- ❖ **Statistical independence:**

$$p(x, y) = p_x(x)p_y(y)$$

- ❖ **Kullback-Leibler divergence (Distance?) of pdf's**

$$L(p(x), p'(x)) = - \int p(x) \ln \frac{p'(x)}{p(x)} dx$$

Pay attention that $L(p(x), p'(x)) \neq L(p'(x), p(x))$

3

- ❖ **Characteristic function** of a pdf:

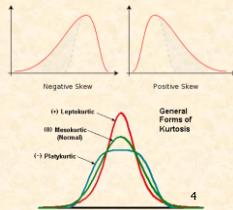
$$\Phi(\Omega) = \int_{-\infty}^{\infty} p(\mathbf{x}) \exp(j\Omega^T \mathbf{x}) d\mathbf{x} = E[\exp(j\Omega^T \mathbf{x})]$$

$$\Phi(s) = \int_{-\infty}^{\infty} p(x) \exp(sx) dx = E[\exp(sx)]$$

- ❖ **2nd Characteristic function:** $\Psi(s) = \ln \Phi(s)$

- ❖ **n-th order moment:** $\frac{d^n \Phi(0)}{ds^n} = E[x^n]$

- ❖ **Cumulants:** $\kappa_n = \frac{d^n \Psi(0)}{ds^n}$



When $E[x] = 0$, then

$$\kappa_0 = 0, \kappa_1 = E[x] = 0,$$

$$\kappa_2 = E[x^2] = \sigma^2, \kappa_3 = E[x^3] \text{ (Skewness)}$$

$$\kappa_4 = E[x^4] - 3\sigma^4 \text{ (Kurtosis)}$$

Discrete Distributions

- ❖ **Binomial distribution B(n,p):**

Repeatedly grab n balls, each with a probability p of getting a black ball. The probability of getting k black balls:

$$P(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

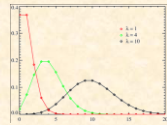
- ❖ **Poisson distribution**

probability of # of events occurring in a fixed period of time if these events occur with a known average.

$$P(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$$

When $n \rightarrow \infty$ and np remains constant,

$$B(n, p) \rightarrow \text{Poisson}(np)$$



Normal (Gaussian) Distribution

- ❖ **Univariate N(μ, σ²):**

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

- ❖ **Multivariate N(μ, Σ):**

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

with the mean μ and the covariance matrix

$$\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12} & \dots & \sigma_{1j} \\ \sigma_{21} & \sigma_{22}^2 & \dots & \sigma_{2j} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{j1} & \sigma_{j2} & \dots & \sigma_{jj}^2 \end{bmatrix}$$

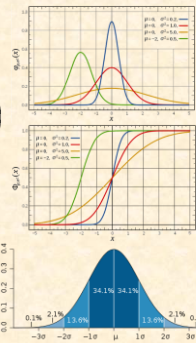
where $\sigma_i^2 = E[(x_i - \mu_i)^2]$ and

$$\sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)]$$

- ❖ **Central limit theorem:**

Let $z = \frac{\sum_{i=1}^n x_i - n\mu}{\sigma\sqrt{n}} \sim N(0,1)$ when $n \rightarrow \infty$

irrespective of the pdf's of x_i 's.



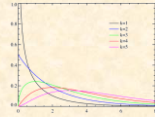
Other Continuous Distributions

- ❖ **Chi-square (χ^2) distribution** of k degrees of freedom: distribution of a sum of squares of k independent standard normal random variables, that is, $\chi^2 = x_1^2 + x_2^2 + \dots + x_k^2$ where $x_i \sim N(0,1)$

$$p(y) = \frac{1}{2^{k/2} \Gamma(k/2)} y^{k/2-1} e^{-y/2} \text{step}(y),$$

where $\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$

- ❖ Mean: k , Variance: $2k$
- ❖ Assume $x \sim \chi^2(k)$
 - Then $(x-k)/\sqrt{2k} \sim N(0,1)$ as $k \rightarrow \infty$ by central limit theorem.
 - Also $\sqrt{2x}$ is approximately normally distributed with mean $\sqrt{2k-1}$ and **unit variance**.



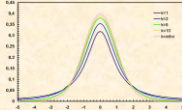
Other Continuous Distributions

- ❖ **t-distribution:** estimating mean of a normal distribution when sample size is small.

A t-distributed variable $q = x/\sqrt{z/k}$ where $x \sim N(0,1)$ and $z \sim \chi^2(k)$

$$p(q) = \frac{\Gamma((k+1)/2)}{\sqrt{k\pi} \Gamma(k/2)} \left(1 + \frac{q^2}{k}\right)^{-(k+1)/2}$$

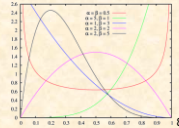
Mean: 0 for $k > 1$,
variance: $k/(k-2)$ for $k > 2$



- ❖ **β -distribution:** Beta(α, β): the posterior distribution of p of a binomial distribution after $\alpha-1$ events with p and $\beta-1$ with $1-p$.

$$p(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$= \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$



Linear Algebra

- ❖ Eigenvalues and eigenvectors: there exists λ and v such that $Av = \lambda v$
- ❖ Real matrix A is called *positive semidefinite* if $x^T A x \geq 0$ for **every** nonzero vector x ;
 A is called *positive definite* if $x^T A x > 0$.
- ❖ Positive definite matrixes act as positive numbers.
All positive eigenvalues
- ❖ If A is symmetric, $A^T = A$,
then its eigenvectors are orthogonal, $v_i^T v_j = 0$.
- ❖ Therefore, a symmetric A can be diagonalized as
 $A = \Phi \Lambda \Phi^T$ and $\Phi^T A \Phi = \Lambda$
where $\Phi = [v_1, v_2, \dots, v_n]$ and $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$

Correlation Matrix and Inner Product Matrix

Principal component analysis (PCA)

- Let x be a random variable in \mathbb{R}^n , its correlation matrix $\Sigma = E[xx^T]$ is positive semidefinite and thus can be diagonalized as

$$\Sigma = \Phi \Lambda \Phi^T$$

- Assign $x' = \Phi^T x$, then $\Sigma' = E(x'x'^T) = \Phi^T \Sigma \Phi = \Lambda$
- Further assign $x'' = \Lambda^{-1/2} \Phi^T x$, then $\Sigma'' = E(x''x''^T) = I$

Classical multidimensional scaling (classical MDS)

- Given a distance matrix $D = \{d_{ij}\}$, the inner product matrix $G = \{x_i^T x_j\}$ can be computed by a bidirectional centering process

$$G = -\frac{1}{2} \left(I - \frac{1}{n} ee^T \right) D \left(I - \frac{1}{n} ee^T \right) \text{ where } e = [1, 1, \dots, 1]^T$$

- G can be diagonalized as $G = \Psi \Lambda' \Psi^T$
- Actually, $n\Lambda$ and Λ' share the same set of eigenvalues, and $\Phi = X^T \Psi$ where $X = [x_1, \dots, x_n]^T$

Because $G = XX^T$, X can then be recovered as $X = \Psi \Lambda'^{1/2}$

Cost Function Optimization

- Find θ so that a differentiable function $J(\theta)$ is minimized.

Gradient descent method

- Starts with an initial estimate $\theta(0)$
- Adjust θ iteratively by

$$\theta_{new} = \theta_{old} + \Delta \theta$$

$$\Delta \theta = -\mu \frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta=\theta_{old}}, \text{ where } \mu > 0$$

- Taylor expansion of $J(\theta)$ at a stationary point θ^0

$$J(\theta) = J(\theta^0) + (\theta - \theta^0)^T \mathbf{g} + \frac{1}{2} (\theta - \theta^0)^T \mathbf{H} (\theta - \theta^0) + O((\theta - \theta^0)^3)$$

$$\text{where } \mathbf{g} = \frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta=\theta^0} \text{ and } \mathbf{H}(i, j) = \frac{\partial^2 J(\theta)}{\partial \theta_i \partial \theta_j} \Big|_{\theta=\theta^0}$$

Ignore higher order terms within a neighborhood of θ^0

$$\theta_{new} - \theta^0 = (I - \mu \mathbf{H})(\theta_{old} - \theta^0)$$

\mathbf{H} is positive semidefinite, then $\mathbf{H} = \Phi \Lambda \Phi^T$, we get

$$\Phi^T (\theta_{new} - \theta^0) = (I - \mu \Lambda) \Phi^T (\theta_{old} - \theta^0)$$

which will converge if every $|1 - \mu \lambda_i| < 1$, i.e., $\mu < 2 / \lambda_{max}$.

Therefore, the convergence speed is decided by $\lambda_{min} / \lambda_{max}$.



Newton's method

- Adjust θ iteratively by

$$\Delta \theta = -\mathbf{H}^{-1} \frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta=\theta_{old}}$$

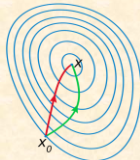
- Converges much faster than gradient descent.

In fact, from the Taylor expansion, we have

$$\frac{\partial J(\theta)}{\partial \theta} = \mathbf{H}(\theta - \theta^0)$$

$$\theta_{new} = \theta_{old} - \mathbf{H}^{-1}(\mathbf{H}(\theta_{old} - \theta^0)) = \theta^0$$

- The minimum is found in one iteration.

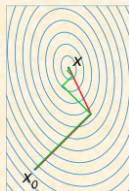


Conjugate gradient method

$$\Delta \theta_i = g_i - \beta_i \Delta \theta_{i-1}$$

$$\text{where } g_i = \frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta=\theta_i}$$

$$\text{and } \beta_i = \frac{g_i^T g_i}{g_{i-1}^T g_{i-1}} \text{ or } \beta_i = \frac{g_i^T (g_i - g_{i-1})}{g_{i-1}^T (g_i - g_{i-1})}$$



12

Constrained Optimization with Equality Constraints

Minimize $J(\theta)$
subject to $f_i(\theta)=0$ for $i=1, 2, \dots, m$

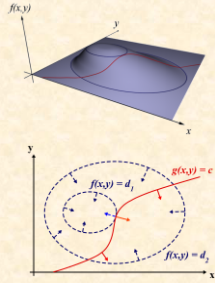
❖ Minimization happens at

$$\frac{\partial J(\theta)}{\partial \theta} = \lambda \frac{\partial f_i(\theta)}{\partial \theta}$$

❖ **Lagrange multipliers:** construct

$$L(\theta, \lambda) = J(\theta) - \sum_{i=1}^m \lambda_i f_i(\theta)$$

and solve $\frac{\partial L(\theta, \lambda)}{\partial \theta} = \frac{\partial L(\theta, \lambda)}{\partial \lambda} = 0$



13

Constrained Optimization with Inequality Constraints

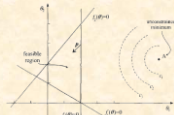
Minimize $J(\theta)$ subject to $f_i(\theta) \geq 0$ for $i=1, 2, \dots, m$

❖ $f_i(\theta) \geq 0$ $i=1, 2, \dots, m$ defines a feasible region in which the answer lies.

❖ **Karush-Kuhn-Tucker (KKT) conditions:**

A set of necessary conditions, which a local optimizer θ_* has to satisfy. There exists a vector λ of Lagrange multipliers such that

- (1) $\frac{\partial}{\partial \theta} L(\theta_*, \lambda) = 0$
- (2) $\lambda_i \geq 0$ for $i=1, 2, \dots, m$
- (3) $\lambda_i f_i(\theta_*) = 0$ for $i=1, 2, \dots, m$



- (1) Most natural condition.
- (2) $f_i(\theta_*)$ is inactive if $\lambda_i = 0$.
- (3) $\lambda_i \geq 0$ if the minimum is on $f_i(\theta_*)$.
- (4) The (unconstrained) minimum in the interior region if all $\lambda_i = 0$.
- (5) For convex $J(\theta)$ and the region, local minimum is global minimum.
- (6) Still difficult to compute. Assume some $f_i(\theta_*)$'s active, check $\lambda_i \geq 0$.

14

❖ **Convex function:**

$f(\theta) : S \subseteq \mathbb{R}^l \rightarrow \mathbb{R}$ is convex if $\forall \theta, \theta' \in S, \lambda \in [0,1]$

$$f(\lambda\theta + (1-\lambda)\theta') \leq \lambda f(\theta) + (1-\lambda)f(\theta')$$

❖ **Concave function:**

$$f(\lambda\theta + (1-\lambda)\theta') \geq \lambda f(\theta) + (1-\lambda)f(\theta')$$

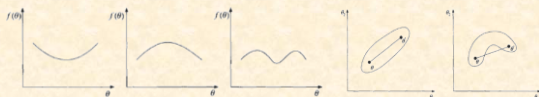
❖ **Convex set:**

$S \subseteq \mathbb{R}^l$ is a convex set if $\forall \theta, \theta' \in S, \lambda \in [0,1]$

$$\lambda\theta + (1-\lambda)\theta' \in S$$

Local minimum of a convex function is also global minimum.

If $f(\theta)$ is concave, then $X = \{\theta \mid f(\theta) \geq b\}$ is a convex set.



15

❖ **Min-Max duality**

Game : A pays $F(x, y)$ \$ to B while A chooses x and B chooses y
 A's goal : $\min_x \max_y F(x, y)$, B's goal : $\max_y \min_x F(x, y)$

The two problems are dual to each other.

In general : $\min_x F(x, y) \leq F(x, y) \leq \max_y F(x, y)$

Therefore, $\max_y \min_x F(x, y) \leq \min_x \max_y F(x, y)$

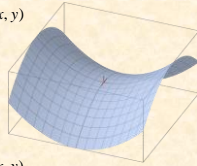
Saddle point condition :

If there exists (x_*, y_*) such that

$$F(x_*, y) \leq F(x_*, y_*) \leq F(x, y_*)$$

or equivalently :

$$F(x_*, y_*) = \max_y \min_x F(x, y) = \min_x \max_y F(x, y)$$



❖ **Lagrange duality**

➤ Recall the optimization problem:

$$\text{Minimize } J(\theta) \quad \text{s.t. } f_i(\theta) \geq 0 \text{ for } i=1,2,\dots,m$$

$$\text{Lagrange function : } L(\theta, \lambda) = J(\theta) - \sum_{i=1}^m \lambda_i f_i(\theta)$$

Because $\max_{\lambda \geq 0} L(\theta, \lambda) = J(\theta)$, we have

$$\min_{\theta} J(\theta) = \min_{\theta} \max_{\lambda \geq 0} L(\theta, \lambda)$$

➤ Convex Programming

For a large class of applications, $J(\theta)$ is convex, $f_i(\theta)$'s are concave then, the minimization on solution (θ, λ) is a saddle point of $L(\theta, \lambda)$

$$L(\theta, \lambda) \leq L(\theta, \lambda) \leq L(\theta, \lambda)$$

$$L(\theta, \lambda) = \min_{\theta} \max_{\lambda \geq 0} L(\theta, \lambda) = \max_{\lambda \geq 0} \min_{\theta} L(\theta, \lambda)$$

Therefore, the optimization problem becomes $\max_{\lambda \geq 0} \min_{\theta} L(\theta, \lambda)$, or

$$\max_{\lambda \geq 0} L(\theta, \lambda) \quad \text{subject to } \frac{\partial}{\partial \theta} L(\theta, \lambda) = 0$$

MUCH SIMPLER!

Mercer's Theorem and the Kernel Method

❖ Mercer's theorem:

Let $x \in \mathfrak{H}'$ and given a mapping $\phi(x) \in H$,

(H denotes Hilbert space, i.e. finite or infinite Euclidean space)

the inner product $\langle \phi(x), \phi(y) \rangle$ can be expressed as a **kernel function**

$$\langle \phi(x), \phi(y) \rangle = K(x, y)$$

where $K(x, y)$ is symmetric, continuous, and positive semi-definite.

The opposite is also true.

The kernel method can transform any algorithm that solely depends on the dot product between two vectors to a kernelized version, by replacing dot product with the kernel function. The kernelized version is equivalent to the algorithm operating in the range space of ϕ . Because kernels are used, however, ϕ is never explicitly computed.
