

语音及语言信息处理国家工程实验室

Math Background (II)







中国科学技术大学 安徽科大讯飞信息科技 股份有限公司



Math Review



- Probability & Statistics
 - Bayes' theorem
 - Random variables: discrete vs. continuous
 - Probability distribution: PDF and CDF
 - Statistics: mean, variance, moment
 - Parameter estimation: MLE
- Information Theory
 - Entropy, mutual information, information channel, KL divergence
- Function Optimization
 - Constrained/unconstrained optimization
- Linear Algebra
 - Matrix manipulation



Information Theory

- Claude E. Shannon (1916-2001, from Bell Labs to MIT): Father of Information Theory, Modern Communication Theory ...
 - C. E. Shannon, "A Mathematical Theory of Communication", Parts 1 & 2, *Bell System Technical Journal*, 1948.
- Information of an event

$$I(A) = \log_2 1/\Pr(A) = -\log_2 \Pr(A)$$

- Entropy(Self-Information) : in bit/nat, amount of info in a R.V.
 - Entropy represents average amount of information in a R.V., the average uncertainty related to a R.V.

语音及语言信息处理国家工程实验室

$$H(X) = -\sum_{x} p(x) \log_2 p(x) = E[\log_2 \frac{1}{p(X)}]$$





Joint and Conditional Entropy

 Joint entropy: average amount of information (uncertainty) about two R.V.s.

$$H(X,Y) = \mathbb{E}[\log_2 \frac{1}{p(X,Y)}] = -\sum_x \sum_y p(x,y) \log_2 p(x,y)$$

 Conditional entropy: average amount of information (uncertainty) of Y after X is known.

$$H(Y \mid X) = \mathbb{E}[\log_2 \frac{1}{p(Y \mid X)}] = -\sum_x \sum_y p(x, y) \log_2 p(y \mid x)$$

• Chain rule for entropy

 $H(X_1, X_2, ..., X_n) = H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, ..., X_{n-1})$

Independence

H(X,Y) = H(X) + H(Y) or H(Y | X) = H(Y)



Mutual Information

• Average information about Y (or X) we can get from X (or Y).

I(X,Y) = H(X) - H(X | Y) = H(Y) - H(Y | X) = H(X) + H(Y) - H(X,Y)

$$I(X,Y) = \sum_{x \in X} \sum_{y \in Y} p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)}$$

or
$$\iint_{x \mid y} p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)} dxdy$$

H(X)
H(X)
H(Y)

• Maximizing I(X,Y) is equivalent to establishing a closer relationship between X and Y, i.e., obtaining a low-noise information channel between X and Y.





• Shannon' s noisy channel model

 $\xrightarrow{\text{Intended}}_{\text{Message}} \xrightarrow{x} \xrightarrow{y} \xrightarrow{p(y/x)} \xrightarrow{p(y$

• A binary symmetric noisy channel



 Channel capacity: the tightest upper bound on the rate of information that can be reliably transmitted over a communication channel

$$C = \max_{p(X)} I(X, Y) = \max_{p(X)} [H(Y) - H(Y | X)] = 1 - H(p) \le 1$$



Example



- X,Y is equiprobable: Pr(X=0)=Pr(X=1)= Pr(Y=0)=Pr(Y=1)= 0.5
- p=0 (noiseless)

$$C = 1 - H(0) = 1$$

• p=0.1(weak noise)

C = 1 - H(0.1) = 0.533

• p=0.4(strong noise)

C = 1 - H(0.4) = 0.03





Bayes' Theorem Application



$$I^* = \arg\max_{I} P(I \mid \hat{O}) = \arg\max_{I} \frac{P(\hat{O} \mid I)P(I)}{P(\hat{O})} = \arg\max_{I} P(\hat{O} \mid I)P(I)$$

Application	Input	Output	p(l)	p(O I)
Speech Recognition	Word Sequence	Speech Features	Language Model (LM)	Acoustic Model
Character Recognition	Actual Letters	Letter images	Letter LM	OCR Error Model
Machine Translation	Source Sentence	Target Sentence	Source LM	Translation (Alignment) Model
Text Understanding	Semantic Concept	Word Sequence	Concept LM	Semantic Model
Part-of-Speech Tagging	POS Tag Sequence	Word Sequence	POS Tag LM	Tagging Model



Kullback-Leibler Divergence



- Introduced by Solomon Kullback and Richard Leibler in 1951
- Distance measure between two PMFs or PDFs (relative entropy)
- D(p||q) ≥0, D(p||q)=0 if and only if q=p, D(p||q) \neq D(q||p)

$$D(p || q) = \mathbb{E}_{p}[\log_{2} \frac{p(x)}{q(x)}] = \sum_{x \in \mathcal{X}} p(x) \log_{2} \frac{p(x)}{q(x)} \text{ or } \int_{x} p(x) \log_{2} \frac{p(x)}{q(x)} dx$$

• Mutual information is a measure of independence

 $I(X,Y) = D(p(x, y) \parallel p(x)p(y))$



Classification: Decision Trees

Example: fruits classification based on features



Classification and Regression Tree (CART)

- Binary tree for classification
 - Each node is attached a YES/NO question
 - Traverse the tree based on the answers to questions
 - Each leaf node represents a class
- CART: Automatically grow a classification tree on a data-driven basis
 - Prepare a finite set of all possible questions
 - For each node, choose the best question to split the node
 - Maximum entropy reduction
 - Small entropy \rightarrow more homogeneous the data is



The CART algorithm

- 1) Question set: create a set of all possible YES/NO questions.
- 2) Initialization: initialize a tree with only one node which consists of all available training samples.
- 3) Splitting nodes: for each node in the tree, find the best splitting question which gives the greatest entropy reduction.
- 4) Go to step 3) to recursively split all its children nodes unless it meets certain stop criterion, e.g., entropy reduction is below a pre-set threshold OR data in the node is already too little.

CART method is widely used in machine learning and data mining:

- 1. Handle categorical data in data mining;
- 2. Acoustic modeling (allophone modeling) in speech recognition;
- 3. Letter-to-sound conversion;
- 4. Automatic rule generation
- 5. etc.



Optimization of objective function (I)

- Optimization:
 - Set up an objective function Q;
 - Maximize or minimize the objective function
- Maximization (minimization) of a function:
 - Differential calculus;
 - Unconstrained maximization/minimization

$$Q = f(x) \vartriangleright \frac{\mathrm{d} f(x)}{\mathrm{d} x} = 0 \vartriangleright x = ?$$
$$Q = f(x_1, x_2, \dots, x_N) \vartriangleright \frac{\P f(x_1, x_2, \dots, x_N)}{\P x_i} = 0 \vartriangleright ?$$

- Lagrange optimization:
 - Constrained maximization/minimization $Q = f(x_1, x_2, \dots, x_N) \text{ with constraint } g(x_1, x_2, \dots, x_N) = 0$ $Q' = f(x_1, x_2, \dots, x_N) + / \times g(x_1, x_2, \dots, x_N)$ $\frac{\P Q'}{\P x_1} = 0, \frac{\P Q'}{\P x_2} = 0, \dots, \frac{\P Q'}{\P x_N} = 0, \frac{\P Q'}{\P x_1} = 0$



Karush–Kuhn–Tucker (KKT) Conditions

• Primary Problem:

 $\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{subject to} \\ g_i(\mathbf{x}) \notin 0 & (i = 1, \cdots, m) \\ h_i(\mathbf{x}) = 0 & (j = 1, \cdots, n) \end{array}$

- Introduce KKT multipliers:
 - For each inequality constraint: M_i $(i = 1, \dots, m)$
 - For each equality constraint: $j (j = 1, \dots, n)$



Karush–Kuhn–Tucker (KKT) Conditions

- Dual problem:
 - if x* is local optimum of the primary problem, x* satisfies:

$$\nabla f(\mathbf{x}^{*}) + \sum_{i=1}^{m} m_{i} \nabla g_{i}(\mathbf{x}^{*}) + \sum_{j=1}^{n} / \sqrt{p} h_{j}(\mathbf{x}^{*}) = 0$$
$$m_{i} \ge 0 \quad (i = 1, \dots, m)$$
$$m_{i} g_{i}(\mathbf{x}^{*}) = 0 \quad (i = 1, \dots, m)$$

• The primary problem can be alternatively solved by the above equations.



Gradient descent (ascent) method:

Optimization of objective function (II)

$$Q = f(x_1, x_2, \cdots, x_N)$$

For any x_i , start from any initial value $x_i^{(0)}$

$$x_i^{(n+1)} = x_i^{(n)} \pm e \cdot \nabla_{x_i} f(x_1, x_2, \cdots, x_N) |_{x_i = x_i^{(n)}}$$

where
$$\nabla_{x_i} f(x_1, x_2, \dots, x_N) = \frac{\partial f(x_1, x_2, \dots, x_N)}{\partial x_i}$$

- Step size is hard to determine
- Slow convergence
- Stochastic gradient descent (SGD)





Optimization of objective function (II)

• Newton' s method

$$Q = f(\mathbf{x})$$

Given any initial value \mathbf{x}_0

$$f(\mathbf{x}) \approx f(\mathbf{x}_{0}) + \nabla f(\mathbf{x}_{0})(\mathbf{x} - \mathbf{x}_{0})^{t} + \frac{1}{2}(\mathbf{x} - \mathbf{x}_{0})^{t} H(\mathbf{x} - \mathbf{x}_{0})$$

$$H = \begin{bmatrix} \frac{\partial^{2} f(\mathbf{x})}{\partial x_{1}^{2}} & \frac{\partial^{2} f(\mathbf{x})}{\partial x_{1} \partial x_{2}} & \cdots & \frac{\partial^{2} f(\mathbf{x})}{\partial x_{1} \partial x_{N}} \\ \frac{\partial^{2} f(\mathbf{x})}{\partial x_{1} \partial x_{2}} & \frac{\partial^{2} f(\mathbf{x})}{\partial x_{2}^{2}} & \cdots & \frac{\partial^{2} f(\mathbf{x})}{\partial x_{2} \partial x_{N}} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial^{2} f(\mathbf{x})}{\partial x_{1} \partial x_{N}} & \frac{\partial^{2} f(\mathbf{x})}{\partial x_{2} \partial x_{N}} & \cdots & \frac{\partial^{2} f(\mathbf{x})}{\partial x_{N}^{2}} \end{bmatrix}_{\mathbf{x} = \mathbf{x}_{0}}$$

$$\mathbf{x}^{*} = \mathbf{x}_{0} - H^{-1} \cdot \nabla f(\mathbf{x}_{0})$$



- Hessian matrix is too big; hard to estimate
- Quasi-Newton' s method: no need to compute Hessian matrix; quick update to approximate it.



Optimization Methods

- Convex optimization algorithms:
 - Linear Programming
 - Quadratic programming (nonlinear optimization)
 - Semi-definite Programming
- EM (Expectation-Maximization) algorithm.
- Growth-Transformation method.





Other Relevant Topics

- Statistical Hypothesis Testing

 Likelihood ratio testing
- Linear Algebra:
 - Vector, Matrix;
 - Determinant and matrix inversion;
 - Eigen-value and eigen-vector
 - Derivatives of matrices;
 - etc.



