

数理解析研究所共同研究会
「凝縮力学系の分子凝縮生物学への応用数理」
期日: 2009年1月5日(月)~6日(金)
場所: 京都大学理学部3号館110号講義室

タンパク質構造と進化と情報幾何

(数学サイドから)

1

- The 1st MSJ-SI: Probabilistic Approach to Geometry (2008.7.28-8.8, 京都大学) での招待講演

“Information divergence geometry and its application to machine learning”
が元になっている。

自己紹介

えぐち しんどう
江口 真透


統計数理研究所・総研大統計科学
eguchi@ism.ac.jp
http://www.ism.ac.jp/~eguchi/

- 数学科出身(大阪大理学部数学科1979年卒)も数学のレールからはずれ、統計学を学ぶ。
- 統計と幾何の融合を目指す。(情報幾何)
- 「忘れ去られた科学: 数学」に関心を持つ。
- 統計教育の見直しの緊急性を感じる。


2

- 大学院から統計の勉強をする。
- 甘利俊一氏が情報幾何を始めた時期で、プレプリントを見て興味を持ち修士のテーマを変える。
- 社会的に見て数学は他分野に比べ進展していない。
- 統計学科は日本にない。したがって日本で統計学を学部で専門的に勉強することは難しい。
- 日本の国家戦略として統計学が抜けているのはマイナスではないか。少なくとも統計学の教育は大事にすべき。

紹介



1. 情報ダイバージェンスの幾何
 - 情報幾何の紹介
 - 情報ダイバージェンスとは
 - 情報ダイバージェンスの連想する幾何
 - 情報ダイバージェンス最小化原理
2. 統計機械学習
 - (カーネル PCA)
 - フースティングと遺伝子発現



- カーネル PCA の話は時間の都合上省略。



目的



3. タンパク質構造と進化との関連

- KL ダイバージェンスと Tsallis ダイバージェンス
- 情報エントロピーと最大分布モデル
- 小山氏の4つの問題提起
 - 構造変化を扱える分子デザイン
 - 構造揺らぎと進化ダイナミクス
 - 分子機能の不可逆性(順序性)
 - 酵素と量子情報幾何

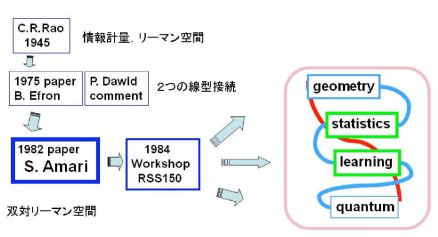




目的は数学と生物の融合

- 小山氏の講演にあった KL ダイバージェンスの最大化の問題を情報幾何の観点から見ていきたい。
- 小山氏による4つの問題は情報幾何の方法論のバラエティの中で扱えるのではないか。

情報幾何の歴史



Timeline of Information Geometry history:

- 1945: C.R. Rao (情報計量, リーマン空間)
- 1975: B. Efron (1975 paper), P. Dawid (comment)
- 1982: S. Amari (1982 paper)
- 1984: RSS150 Workshop

Key concepts: 2つの線形接続, 双対リーマン空間

Central themes: geometry, statistics, learning, quantum

- 歴史的に見て統計学の中での最初の幾何的な考え方は Gauss の最小二乗法, 線形射影という話から Gauss norm の定義が成立する。
- 微分幾何的な考え方は C.R. Rao が最初ではないか (Bull. Calcutta Math. Soc., 1945). 情報計量を元にして統計モデルはリーマン空間になる。
- 実はリーマン空間では不十分。2つの線形接続が必要になる (Efron, Dawid)。
- 甘利氏により標準的な微分幾何の言葉で定式化。

What is IG?

Geometry

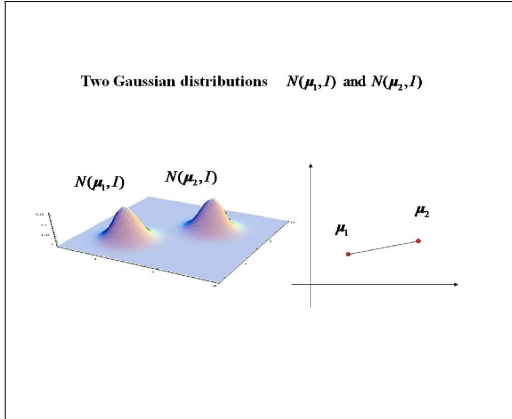
↔

Uncertainty

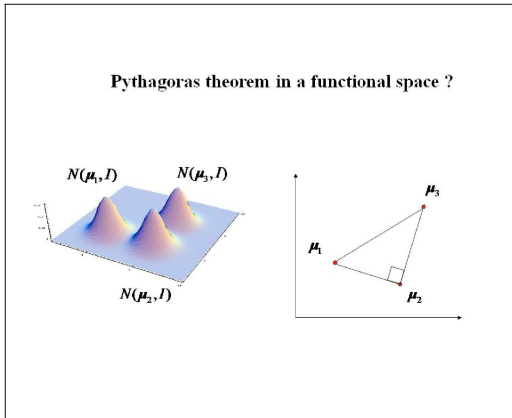
Information space

It is a **method** to quantify uncertainty,
Or, a **viewpoint** to understand uncertainty?

- 情報幾何は情報空間の中での不確定性を幾何的に sharp に理解することを目指している。



- 統計モデルは分布の集まりである。規格化することにより分布をユークリッド空間内の点と思うと、2つの分布を結ぶ path は2点を結ぶ path に相当する。



- 3つの分布があったとき、対応するユークリッド空間内の3点の関係、たとえばピタゴラスの定理が考えられる。

2種類のワンパラメータ族

データ空間の確率密度の全空間を \mathbb{P} とする。

m-測地線 $p_t^{(m)}(x) = (1-t)p(x) + tq(x), \quad (p, q \in \mathbb{P})$

e-測地線 $r_s^{(e)}(x) = c_s \{r(x)\}^{1-s} \{q(x)\}^s, \quad (q, r \in \mathbb{P})$

ここで $c_s = 1 / \int \{r(x)\}^{1-s} \{q(x)\}^s dx$

上記の一般化

- 3点 p, q, r に対し、p と q を結ぶ測地線および r と q を結ぶ測地線を結び、ピタゴラスの定理を考えてみる。

自己紹介

えぐち しんとう
江口 真透

統計数理研究所・総研大統計科学
eguchi@ism.ac.jp
http://www.ism.ac.jp/~eguchi/

- 数学科出身(大阪大理学部数学科1979年卒)も数学のレールからはずれ、統計学を学ぶ。
- 統計と幾何の融合を目指す。(情報幾何)
- 「忘れ去られた科学:数学」に関心を持つ。
- 統計教育の見直しの緊急性を感じる。

2

Kullback-Leibler(KL) ダイバージェンス

- 正規分布に対し、KL ダイバージェンスはユークリッド距離の2乗に比例するのでピタゴラスの定理と同じ意味になる。
- すなわち、3点を同一平面上にとるユークリッド幾何に帰着される。

Pythagorean Theorem

$$\dot{p}_i^{(m)} \perp_{g_i} \dot{r}_i^{(m)} \Rightarrow D_{KL}(p_i^{(m)}, r_i^{(m)}) = D_{KL}(p_i^{(m)}, q_i^{(m)}) + D_{KL}(q_i^{(m)}, r_i^{(m)})$$

$D_{KL}(\dot{p}_i, \dot{r}_i) \geq D_{KL}(q_i, \dot{r}_i)$
 $D_{KL}(\dot{p}_i, \dot{r}_i) \geq D_{KL}(\dot{p}_i, q_i)$ (em algorithm, Amari, 1995)

13

- 2種類の測地線に対し、ピタゴラスの定理が成立するための必要十分条件は2つの path の点 q での接ベクトルが直交していることである。
- これにより2つの不等式、いわゆる em アルゴリズムが簡単に導かれる。

Information Geometry

Statistical model $M = \{q(x, \theta) : \theta \in \Theta\}$
 where $q(x, \theta)$ is a pdf st $\int q(x, \theta) dx = 1$

Information metric (Fisher, 1922; Rao, 1945)
 $g_{ij}(\theta) = E_{q(x, \theta)} \{e_i(x, \theta) e_j(x, \theta)\}$ where $e_i(x, \theta) = \frac{\partial}{\partial \theta^i} \log q(x, \theta)$

Dual connections (Amari, 1982)

Exponential connection $\overset{\circ}{\Gamma}_{ij\lambda}(\theta) = E_{q(x, \theta)} \left(\frac{\partial^2}{\partial \theta^i \partial \theta^j} \log q(x, \theta) e_\lambda(x, \theta) \right)$

Mixture connection $\overset{\circ}{\Gamma}_{ij\lambda}(\theta) = \int \frac{\partial^2}{\partial \theta^i \partial \theta^j} q(x, \theta) e_\lambda(x, \theta) dx$

$\frac{1}{2}(\overset{\circ}{\Gamma}_{ij\lambda}(\theta) + \overset{\circ}{\Gamma}_{ij\lambda}(\theta))$ is Riemannian connection with $\{g_{ij}(\theta)\}$

9

- 統計モデル M を θ でパラメトライズされた分布の集まりとする。
- 情報計量について
 - Fisher がフィッシャー情報量を定義。
 - Rao が³, これを計量とみなすことによりリーマン空間となることを指摘。
 - この量はハイゼンベルグの不確定性原理とほとんど等価で推定の限界を与える。
- 2つの接続 (e-接続, m-接続) が古典的な係数で定義される。
 - 統計モデルの対数密度の二階微分と接ベクトルの内積
 - 密度の二階微分
- 情報幾何はこの2つの接続の組で一つの測地性を与える。
- この2つの接続の平均をとることにより g のもとでのリーマン接続になる。

Mixture and exponential models

Mixture model (mixture geodesic space)

$$M = \{g(x, \theta) : \theta \in \Theta\}$$

$$g(x, \theta) = \sum_{i=1}^k \theta_i g_i(x), \quad \Theta = \{(\theta_1, \dots, \theta_k) : \theta_i > 0, \sum_{i=1}^k \theta_i = 1\}$$

Cf. Maximum likelihood

Exponential model (exponential geodesic space)

$$M = \{g(x, \theta) : \theta \in \Theta\}$$

$$g(x, \theta) = \exp\left\{\sum_{i=1}^k \theta_i f_i(x) - \psi(\theta)\right\}, \quad \Theta = \{\theta : \psi(\theta) < \infty\}$$

where $\psi(\theta) = \log \int \exp\{\theta f_i(x)\} dx$

Cf. Gaussian, Bernoulli, Poisson, Gamma, ... 10

2つの接続に関する測地線・測地面

- どちらも分布・は2種類の接続に対するアフィンパラメータになっている。
- 統計学とは統計モデルの上で統計的推定と統計的検定をする。モデルを作るということと、それで推論する、ということは別。したがって2種類の接続が必要。
 - モデル：対象
 - 推論：行為
- 対象と行為が組になって一つの統計という作業ができる。よって2種類の測地性が必要。
- これはアインシュタインの相対論（重力場の幾何学）と大きく異なる点。

Information divergence geometry

$D : M \times M \rightarrow \mathbf{R}$ is an **information divergence** on a statistical model M

\Leftrightarrow (i) $D(p, q) \geq 0$ with equality if and only if $p = q$
 (ii) D is differentiable on $M \times M$

Let $D(Z \cdots | XY \cdots | p) = Z_p \cdots X_q Y_q \cdots D(p, q)|_{p=p}$

Then we get a **Riemannian metric** and **dual connections** on M (Eguchi, 1983, 1992)

$$g^p, \quad g^p(X, Y) = -D(X|Y)$$

$$\nabla^p, \quad g^p(\nabla^p X, Z) = -D(XY|Z) \quad (\forall Z \in X(M))$$

$${}^* \nabla^p, \quad g^p({}^* \nabla^p Y, Z) = -D(Z|XY) \quad (\forall Z \in X(M))$$

14

ダイバージェンスの定義

- KL ダイバージェンス以外のダイバージェンスも統一的に見ることができないか。
 - 特に対角部分
 - $\bullet = \{(p, p) : p \in M\}$
- で可微分であることが必要。
- ダイバージェンスはリーマン計量 g , 2つの接続 $\bullet, {}^* \bullet$ を定める
 - 微分 $D(Z \cdots | XY \cdots | p)$ は Z と XY について非対称なので、 \bullet と ${}^* \bullet$ は一般に異なる。
 - KL ダイバージェンスの場合は
 - \bullet : e-接続
 - ${}^* \bullet$: m-接続
- となり、2種類の接続がKLの非対称性から出てくる。

Proof 1

g^p is a Riemannian metric

(i) $g^p(X, Y) = g^p(Y, X)$
 $D(X|Y) = 0$ since $D(p, q) \geq D(p, p) = 0$
 $g^p(X, Y) - g^p(Y, X) = D(XY|Y) - D(YX|X) = D[(X, Y)|Y] = 0$

(ii) g^p is positive-definite
 $g^p(X, Y) = -D(X|Y) = D(XY|Y)$
 because $D(Y|Y) = 0$.

Proof 2

∇^p and ${}^* \nabla^p$ are dual connections

(i) $\nabla^p X Y = f \nabla^p X Y$
 (ii) ${}^* \nabla^p X Y = f \nabla^p X Y + (Xf)Y$

$$g^p(\nabla^p X Y - f \nabla^p X Y, Z) = -D((fX)Y - fXY|Z) = 0 \quad (\forall Z)$$

$$g^p(\nabla^p X(fY) - f \nabla^p X Y - (Xf)Y, Z) = -D(X(fY) - (Xf)Y - fXY|Z) = 0 \quad (\forall Z)$$

(iii) Let $\bar{\nabla} = \frac{1}{2}(\nabla^p + {}^* \nabla^p)$.
 Then $X g^p(Y, Z) = g^p(\bar{\nabla} X, Z) + g^p(X, \bar{\nabla} Z)$

$$X g^p(Y, Z) - g^p(\bar{\nabla} X, Z) - g^p(X, \bar{\nabla} Z) = -XD(Y|Z) + \frac{1}{2}\{XD(Y|Z) + XD(Z|Y)\} = 0$$

16

参考書(微分幾何)



リーマン幾何学(近代数学講座)(単行本)
立花 俊一(著) 単行本:
出版社 朝倉書店, 復刊版 (2004/04)
ISBN-10: 4254116586



現代微分幾何入門(単行本)
野水 克己(著)
単行本:
出版社 養華房 (1981/01)
ISBN-10: 4785311274

Foundations of Differential Geometry (Wiley
Classics Library) (ペーパーバック)
Shoshichi Kobayashi (著), Katsumi Nomizu (著)

リーマン幾何学(近代数学講座)(単行本)
立花 俊一(著) 単行本:
出版社 朝倉書店, 復刊版 (2004/04)
ISBN-10: 4254116586

現代微分幾何入門(単行本)
野水 克己(著)
単行本:
出版社 養華房 (1981/01)
ISBN-10: 4785311274

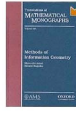
Foundations of Differential Geometry (Wiley
Classics Library) (ペーパーバック)
Shoshichi Kobayashi (著), Katsumi Nomizu (著)

どんな参考書(情報幾何)?



Differential Geometrical Methods in Statistics
Shun-ichi Amari 1985 年 Springer

<http://www.geocities.jp/imitats/eqnchi.pdf/okoo.pdf>



Methods of Information Geometry
Shun-ichi Amari, Hiroshi Nagacka
Amer. Mathematical Society (2001)

Differential Geometrical Methods in Statistics
Shun-ichi Amari 1985 年 Springer

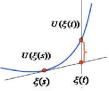
<http://www.geocities.jp/imitats/eqnchi.pdf/okoo.pdf>

Methods of Information Geometry
Shun-ichi Amari, Hiroshi Nagacka
Amer. Mathematical Society (2001)

U divergence

Let us take a convex function $U: \mathbb{R} \rightarrow \mathbb{R}$.
Let $u(z) = \frac{1}{\alpha} U(z)$, $u(\xi(z)) = z$.

$\xi(s) = \arg \max_{-\infty < t < \infty} \{ts - U(t)\}$
 $u(t) = \arg \max_{-\infty < s < \infty} \{st - U^*(s)\}$
where $U^*(s) = s\xi(s) - U(\xi(s))$



U divergence

$$D_U(p, q) = \int \{U(\xi(q)) - U(\xi(p)) - u(\xi(q))(\xi(q) - \xi(p))\}$$

$$= \int \{U(\xi(q)) - p\xi(q) + U^*(\xi(p))\}$$

19

U-ダイバージェンスの構成法

- u は u の逆関数.
- U-ダイバージェンスは s と t として分布 p と q を代入することにより得られる.
- $u(\cdot(p)) = p$ に注意. これは U-ダイバージェンスの持っているきれいな性質.

Example of U divergence

$$D_U(p, q) = \int \{U(\xi(q)) - U(\xi(p)) - p(\xi(q) - \xi(p))\}$$

Kullback-Leibler (KL) divergence
 $U(t) = \exp(t)$

$$D_{KL}(p, q) = \int \{q - p - p(\log q - \log p)\} = \int p \log \frac{p}{q}$$

Tsallis power divergence
 $U_p(t) = \frac{(1 + \beta t)^{1/\beta}}{1 + \beta}$

$$D_\beta(p, q) = \int \left\{ \frac{q^{\beta+1} - p^{\beta+1}}{\beta+1} - p \left(\frac{q^\beta - p^\beta}{\beta} \right) \right\}$$

Note $\lim_{\beta \rightarrow 1} U_p(t) = \exp(t)$
 $\lim_{\beta \rightarrow 0} D_\beta(p, q) = D_{KL}(p, q)$

20

典型的な U-ダイバージェンス

- U を指数関数とすると KL ダイバージェンス.
- U-ダイバージェンスの中でも KL が一番いい. すなわちデータの分布が仮定された分布に exact に従うとすると U が指数であるというのが一番いい.
- Tsallis power ダイバージェンスもべき指数を使って構成できるので U-ダイバージェンスの一種. これはアウトライヤーの重みを 0 にするのでデータが汚れているときに有用.
- 現実においては仮定した統計モデルと現実とは乖離があるので, それにあわせて Tsallis power ダイバージェンスの β を決める, という事が行われている.
- とはいえ, 実際には U の選び方はそんなに効かない. U の種類もそんなに知られていない.

U-エントロピー

U-クロスエントロピー $C_U(p, q) = -E_p\{\xi(q)\} + \langle U(\xi(q)) \rangle$

U-エントロピー $H_U(p) = -E_p\{\xi(p)\} + \langle U(\xi(p)) \rangle$

U-ダイバージェンス $D_U(p, q) = C_U(p, q) - H_U(p)$

例題 1. Let $U(t) = \exp(t)$
 Boltzmann-Shannonエントロピー $H_U(p) = E_p(-\log p)$

例題 2. Let $U(t) = \frac{\alpha + \beta t}{1 + \beta t}$
 Tsallisエントロピー $H_U(p) = E_p\left(-\frac{\beta(X)-1}{\beta}\right)$

ダイバージェンスの分解

- クロスエントロピーとエントロピーに分解される.
- $\langle \cdot \rangle$ は積分を表す.
- E_p は p に関する期待値.

Geometric formula with D_U

$(g^{(U)}, \nabla^{(U)}, \nabla^{*(U)})$ s.t. $g^{(U)}(X, Y) = -D_U(X|Y)$
 $g^{(U)}(\nabla_X Y, Z) = -D_U(XY|Z) \quad (\forall Z \in X(M))$
 $g^{(U)}(\nabla_X^* Y, Z) = -D(Z|XY) \quad (\forall Z \in X(M))$

$g_{ij}^{(U)}(\theta) = \int \frac{\partial}{\partial \theta^i} q(x, \theta) \frac{\partial}{\partial \theta^j} \xi(q(x, \theta)) dx$
 $\Gamma_{jk}^{(U)}(\theta) = \int \frac{\partial^2}{\partial \theta^j \partial \theta^k} q(x, \theta) \frac{\partial}{\partial \theta^i} \xi(q(x, \theta)) dx$
 $\Gamma_{jk}^{*(U)}(\theta) = \int \frac{\partial}{\partial \theta^k} q(x, \theta) \frac{\partial^2}{\partial \theta^i \partial \theta^j} \xi(q(x, \theta)) dx$

$g^{(U)} = g \iff U = \exp$ (KL divergence)
 $\nabla^{(U)} = \nabla^m$ (∇U) (Eguchi, 2005)

22

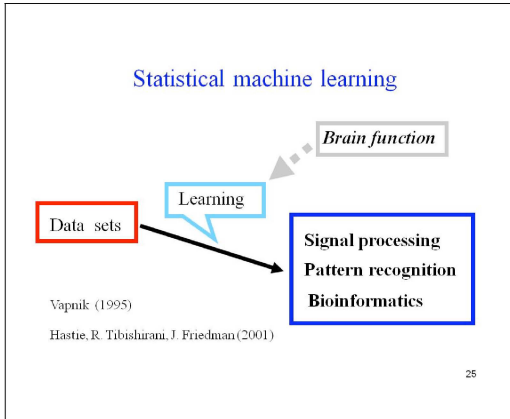
ダイバージェンスが作る幾何量である計量と組の接続を公式に基づいて計算

- パラメータを作って成分で計算.
- (任意の U に対し) $\bullet_{ij,k}^{(U)}$ は m -接続になっている. これは U -ダイバージェンスの特徴であり, 統計的な方法に使うときに便利な性質である.

Application of
 minimum divergence method
 to Machine learning

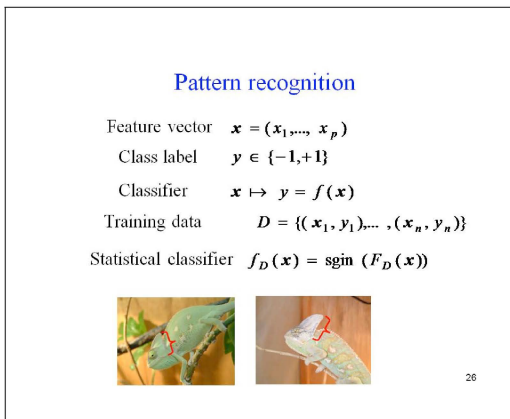
24

統計的な応用について



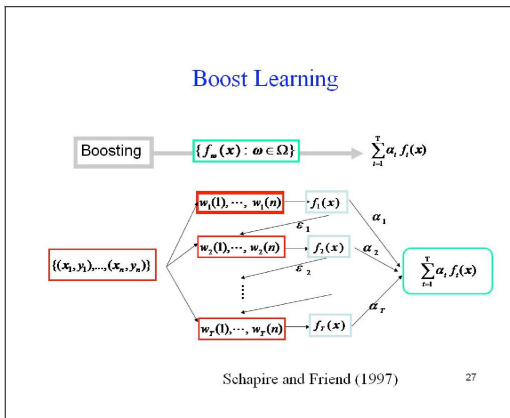
統計機械学習について

- ここ 4・5 年の統計の応用のはやり
- 「データをもとにして統計的な方法を行う」事を 学習する という (それ以前は「推定する」「検定する」).
- 統計機械学習はある意味、脳の機能を imitate して学習する.



パターン認識について

- 識別子 f は \mathbf{x} を見てラベルを予測する.
- 最終的に判別関数 $F_D(\mathbf{x})$ を作って、その符号によりラベルをわけける.
- たとえばオスのエボシカメレオンがメスを獲得できるか、をラベルとすると識別子は頭頂のクレスト.



ブースト学習

- SVM (Support Vector Machine) に対抗.
- 識別子の辞書 $\{f_i(\mathbf{x}) : \dots\}$ を用意し、うまく多数決、すなわち重み α_i を付けて線形和をとることにより、よい統計分類を作る. このとき、辞書をたくさん用意すると、ロスを汎関数的に最小化できる.
- n 個の例題に対し、毎回重みを変えると、一番よかった識別子が毎回異なるので、それらの識別子を統合して多数決をとる.
- ポイントは重み $w_1(1), \dots, w_r(n)$ の付け方.
 1. $w_1(1) = \dots = w_1(n) = 1/n$.
 2. $f_1(\mathbf{x})$ が誤識別した例題のみ非常に大きくなるように重みを決定.
 3. 以降、Markov 的に重みを決定.

U-boost learning

Let (x, y) be a variable with feature vector x and label y .

Decision rule $h(x) = \arg \max_{y \in \{0, -1\}} \hat{q}(y|x)$

U-boost: $\hat{q}_t(y|x) = u(\alpha_t f_t(x, y) + \xi(q_{t-1}(y|x)))$

Step 1. $f_t = \arg \min_f \zeta(f|q_{t-1})$

Step 2. $\alpha_t = \arg \min_{\alpha} L_{\sigma}^{\text{emp}}(\alpha f_t + \xi(q_{t-1}))$

$L_{\sigma}^{\text{emp}}(\xi(q_t)) - L_{\sigma}^{\text{emp}}(\xi(q_{t-1})) = D_{\sigma}(q_{t-1}, q_t)$

Murata et al (2004) 28

U-ダイバージェンスにブースト学習を応用

- σ_t は U-ダイバージェンスが作るロスの gradient.
- $L_{\sigma}^{\text{emp}}(\sigma(q_t)) - L_{\sigma}^{\text{emp}}(\sigma(q_{t+1}))$ はステップ t とステップ $t+1$ でのロスの減少幅. これが U-ダイバージェンスの 2 点間の距離と等しい, というピタゴラスの定理が成り立つ.

Simulation (complete separation)

Feature space
[-1, 1] × [-1, 1]

Decision boundary
 $x_2 = \sin(2\pi x_1)$

$\{(x_i, y_i) : i=1, \dots, 1000\}$
 $x_i \in [-1, 1] \times [-1, 1]$
 $y_i \in \{-1, +1\}$

28

非線形な例：サイン波の上が +1, 下が -1

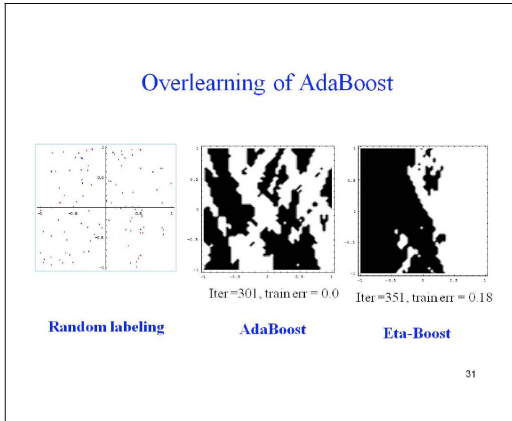
Learning process

Iter = 1, train err = 0.21 Iter = 13, train err = 0.18 Iter = 17, train err = 0.10

Iter = 23, train err = 0.10 Iter = 31, train err = 0.095 Iter = 47, train err = 0.008

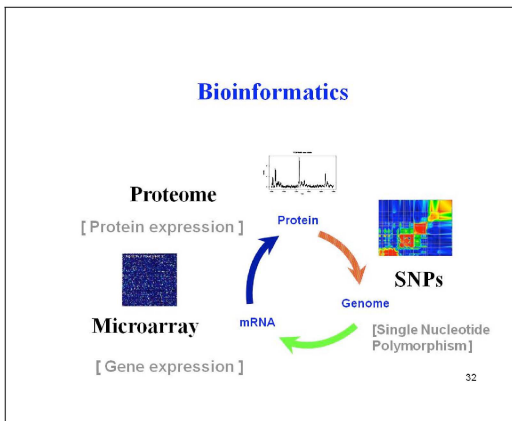
30

- 辞書は線形判別関数のみなので Step1 では線形.
- ブースト学習アルゴリズムが徐々に非線形性を学習する.



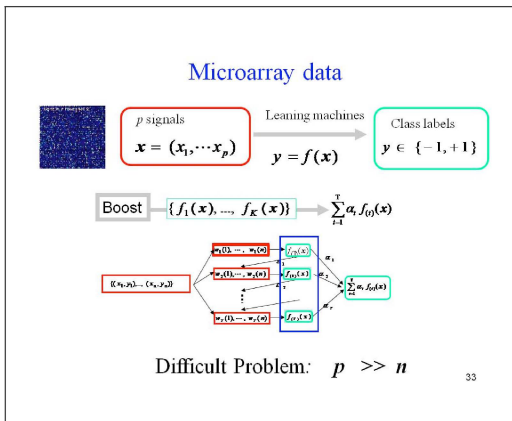
解がない場合、どう学習をするか

- ランダムな点, ランダムなラベル.
- アダブーストでは中央図のようになんとかあわせようとする. 過剰学習が起きる.
- η を少し変形すると右図のようになる. ひとかたまりになるので過剰学習に強い.
- ブーストの収束の判断は難しい.
 - 過剰学習が起きるので, 回数を多くすればいいというものではない.
 - よくやられているのが情報量規準だが, ブースティングでまだよい情報量規準が出ていない.
 - 現在用いている方法:
トレーニングデータを仮想的にトレーニングに使う分と残りにわけてテストする. これをすべての分け方でいい, テストエラーを小さくする回数を見つける.



関連解析

- バイオインフォマティクスの特別なタスク

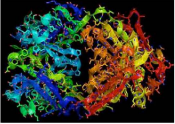


- p 個の遺伝子の発現量.
- y は治療効果があるかどうか (テーラーメイド医療に関連).
- n はサンプルサイズ.
- 医学的に解釈がつくような K を取り, ブーストアルゴリズムにかける.
- p は 2 万~3 万, n は大体 100 以下. 一方, p が 1 万だったら n が 100 万位だと統計的に有意な事ができる. したがってかなり無理なことをしている (p · n 問題). エビデンスを他の外的なインフォメーション (パスウェイの知識, ネットワークの知識など) から持てこないといけない.
- ブーストの辞書の選び方は発現の上位のものから. ただし, 医学的に意味にないのも多いのでジェンリストに限って解析することもある.
- 各 f_i は一つの遺伝子・プローブしか見ていないので, 線形な判別というよりステップワイズな判別.

乳がんのある治療薬の奏功性予測

国立がんセンターの小泉、田村、清水医師の訪問を受ける。
国立がんセンターグループからマイクロアレイの32Gのデータが届く。
プリチャード、小森、江口が国立がんセンターグループに訪問。

トラスツズマブ(ハーセプチン)分子標的治療薬の一種



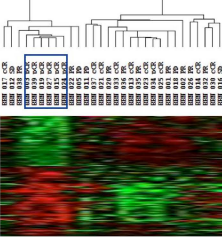
HER2過剰発現が確認された転移性乳癌に対する治療薬

ハーセプチンは、米国において1998年に認可。日本では、2001年4月にHER2過剰発現が確認された転移性乳癌に対する治療薬として承認された。

出典: フリー百科事典『ウィキペディア (Wikipedia)』

- 通常 2~3 年目でデータの半分, 5 年目ですべてのデータが届く。
- ハーセプチンはある乳癌の患者には全く効かない。あらかじめ遺伝子発現量を見て、効くかどうか予測できればテーラーメイド医療につながる。

Clustering analysis



26000 genes

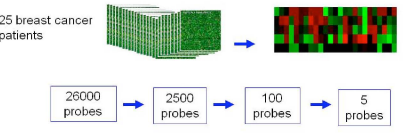
National Cancer Center, Japan

35

階層的クラスタリング

- 青く囲った部分は pCR, すなわち病理学的に完全に治った患者なので, 非常にきれいなデータ (これが vague に出ているとどうしようもない).
- 考慮しなかったが, 本来ノイズにかなり影響されてしまうので実際に医療に使うには, たとえばクラスターをグループとしてもう少し安定な指標にする必要がある.
- アレイのサイエンスが進めば, ラーニングがスーパーバイズド・アンスーパーバイズドを行ったり来たり, ということで悩まなくてもよくなるのでは.

Boost Procedure



25 breast cancer patients

26000 probes → 2500 probes → 100 probes → 5 probes

Boost

5 genes $x = (x_1, \dots, x_5)$

$$f_j(x) = \begin{cases} +1 & \text{if } x_j \leq c_j \\ -1 & \text{if } x_j > c_j \end{cases}$$

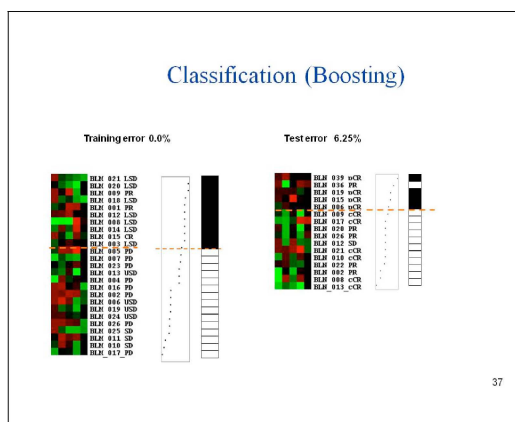
$y = f(x)$

Treatment effect $y \in \{-1, +1\}$

36

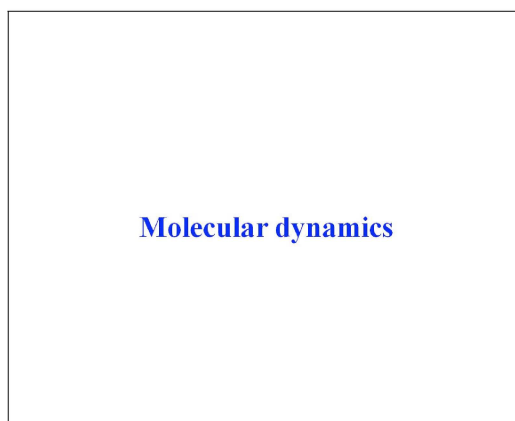
フィルタリング

- 発現が多く, データマイニングして遺伝子の機能がわかっているものを優先的に選択する.
- 最終的に 5 個の遺伝子での予測の成績が一番よかった (6 だと急にパフォーマンスが落ちる). しかし, 必ずしも機能がわかっている本命が出てきていない. 統計的にフィットする予測がいいのか, それともアソシエーションがちゃんとつかまえられるのがいいのか, というのは別の問題.
- 予測に対する多重性の理論はできていないが, 経験的に言ってこれだけうまくいく例はあまりない.
- 26000 個の遺伝子の中で選択した 5 個と似たような遺伝子もあわせて見る, という議論も現在進めている.



トレーニングデータとテストデータ

- アダブーストではいつでもトレーニングエラーを0にできる(ただし、その前にやめた方がいい場合もある)。
- さらに追加のデータに対しては時間差を考慮するなど改良する必要がある。



小山氏のメール・論文に関するコメント

Molecular dynamics

Koyama *et al. Physical Review E* 78, 2008

Perturbed distribution $\rho_2(\mathbf{q}) = \frac{1}{Z} \exp\left(\sum_{i=1}^M (\lambda_i V_i(\mathbf{q}) - \psi(\lambda_i))\right)$

KL-divergence $D_{\text{KL}}(\rho_2, \rho_0) = \langle \lambda \cdot V(\mathbf{q}) - \psi(\lambda) \rangle_2 = \lambda \cdot \mu - \psi(\lambda)$

where $\mu = \langle V(\mathbf{q}) \rangle_2 = \frac{\partial}{\partial \lambda} \psi(\lambda)$

Covariance matrix of $V(\mathbf{q})$ under distribution $\rho_2(\mathbf{q})$

$$C_2 = \langle (V(\mathbf{q}) - \mu)(V(\mathbf{q}) - \mu)^T \rangle_2 = \frac{\partial^2}{\partial \lambda \partial \lambda^T} \psi(\lambda) = \frac{\partial}{\partial \lambda} \mu^T$$

小山氏の論文

Max KL-divergence

Max KL-divergence on e-geodesic sphere

$$\max \{ D_{\text{KL}}(\rho_\lambda, \rho_0) : \lambda \cdot \lambda = 1 \}$$

- 小山氏のしているのは制約付きの KL ダイバージェンスの最大化,
- 一方, 統計で行うのは最小化, それが最尤法と等価, というのがもともとの KL ダイバージェンスの使い方.

PCA proposed by Koyama et al.

Max KL-divergence on e-geodesic sphere

$$\max \{ D_{\text{KL}}(\rho_\lambda, \rho_0) : \lambda \cdot \lambda = 1 \}$$

Approximate ellipsoid to KL

$$D_{\text{KL}}(\rho_\lambda, \rho_0) = \frac{1}{2} \lambda^T C_0 \lambda + O(\|\lambda\|^3)$$

$$\arg \max \{ D_{\text{KL}}(\rho_\lambda, \rho_0) : \lambda \cdot \lambda = 1 \} \approx \arg \max \{ \lambda^T C_0 \lambda : \lambda \cdot \lambda = 1 \}$$

Max KL-divergence on e-geodesic sphere \approx PCA with C_0

- 二次近似して, C_0 を推定して PCA をすることにより KL ダイバージェンスの最大化と近似的に等価,

Iterative PCA?

$$\max \{ D_{\text{KL}}(\rho_\lambda, \rho_0) : \lambda \cdot \lambda = 1 \}$$

Lagrange function $L(\lambda) = D_{\text{KL}}(\rho_\lambda, \rho_0) - \frac{1}{2} \kappa(\lambda \cdot \lambda - 1)$

Gradient $\frac{\partial}{\partial \lambda} L(\lambda) = \frac{\partial}{\partial \lambda} \{ \lambda \cdot \mu - \psi(\lambda) - \frac{1}{2} \kappa(\lambda \cdot \lambda - 1) \} = C_2 \lambda - \kappa \lambda$

$$\lambda^* = \arg \max_{\lambda \cdot \lambda = 1} D_{\text{KL}}(\rho_\lambda, \rho_0) = \arg \max_{\lambda \cdot \lambda = 1} \lambda^T C_2 \lambda$$

Exact equivalence with PCA

Iterational algorithm for finding λ^*

$$\lambda_{i+1} \leftarrow \text{Principal-eigenvector}(\hat{C}_i)$$

$$\hat{C}_i \lambda_{i+1} = \kappa \lambda_{i+1}$$

Exact に計算

- λ^* を求めるために C_2 の PCA
- C_2 と λ^* は入れ子になってしまっている.
 - C_2 を与えると第 1 固有値 λ^* が定義される,
 - C_2 は λ^* の関数,
- 小山氏は摂動が小さいときには摂動前の情報から予測ができるので C_0 としていた,
- 別のアイデアとして, 反復アルゴリズムを用いることにより exact な KL ダイバージェンスの最大化が得られるのではないか,

Boost Learning Algorithm?

Objective functional $L(V) = \langle V(\mathbf{q}) \rangle_{\rho^0} - \psi(V) - \frac{1}{2} \kappa (\|V\|^2 - 1)$

where $\rho^0(\mathbf{q}) = \rho_0(\mathbf{q}) \exp\{\psi(V) - \psi(V)\}$

Dictionary of potential energy functions $F = \{f_\alpha(\mathbf{q}) : \alpha \in A\}$

Boost learning algorithm

$$V_i(\mathbf{q}) \leftarrow V_{i-1}(\mathbf{q}) + \lambda_i f_i(\mathbf{q})$$

$$(i) f_i = \arg \max_{\{f \in F\}} \{ \frac{\partial}{\partial \alpha} L(V_{i-1} + \lambda f) |_{\alpha=0} \}$$

$$(ii) \lambda_i = \arg \max_{\{\lambda \in \mathbb{R}\}} \{ L(V_{i-1} + \lambda f_i) \}$$

$$V(\mathbf{q}) = \sum_{i=1}^T \lambda_i V_i(\mathbf{q})$$

43

ブースト学習が可能

- あまり時間はかかからないのでは、

Max dual KL-divergence?

Max dual KL-divergence on e-geodesic sphere

$$\max \{ D_{KL}(\rho_0, \rho_\lambda) : \lambda \cdot \lambda = 1 \}$$

Lagrange function $L_\lambda(\lambda) = D_{KL}(\rho_0, \rho_\lambda) - \frac{1}{2} \kappa (\lambda \cdot \lambda - 1)$

Gradient $\frac{\partial}{\partial \lambda} L_\lambda(\lambda) = \frac{\partial}{\partial \lambda} \{ -\lambda \cdot \mu_0 + \psi(\lambda) - \frac{1}{2} \kappa (\lambda \cdot \lambda - 1) \} = -\mu_0 + \mu - \kappa \lambda$

$$\mu_\kappa^+ - \kappa^+ \lambda_\kappa^+ = \mu_0, \quad \lambda_\kappa^+ \cdot \lambda_\kappa^+ = 1$$

44

KL の非対称性より \bullet_0 と \bullet_\bullet を逆にする

- PCA の形にならない、
- 少し難しいかもしれないが停留点を見つければ計算できるのでは、
- どれくらい意味があるかはよくわからない、

Max Tsallis power divergence?

Tsallis power divergence

$$D_\beta(\rho_\lambda, \rho_0) = \int \left\{ \frac{\rho_0^{\beta+1} - \rho_\lambda^{\beta+1}}{\beta+1} - \rho_\lambda \frac{\rho_0^\beta - \rho_\lambda^\beta}{\beta} \right\}$$

Tsallis entropy distribution

$$\rho_\lambda(\mathbf{q}) = \frac{1}{Z} \exp_\beta \left(\sum_{i=1}^M (1 + \lambda_i) V_i(\mathbf{q}) - \psi_\beta(\lambda) \right)$$

$$\text{where } \exp_\beta(V) = (1 + \beta V)^{\frac{1}{\beta}}$$

Max Tsallis power divergence on e-geodesic sphere ?

$$\max \{ D_\beta(\rho_\lambda, \rho_0) : \lambda \cdot \lambda = 1 \}$$

45

Tsallis power ダイバージェンスの最大化も考えられる