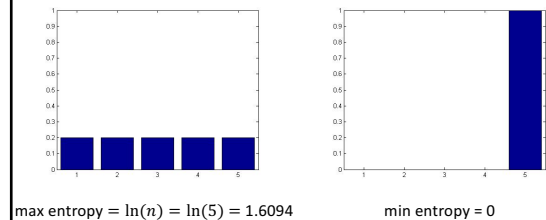


Entropy and related information theoretical measures: Examples on application in data analysis

Allan Aasbjerg Nielsen
Technical University of Denmark
DTU Compute – Department of Applied
Mathematics and Computer Science
alan@dtu.dk

IACG Christmas WS 17 Dec 2013

Entropy



Entropy

- Discrete stochastic variable X with probability density function (pdf) $p(X = x_i)$, $i = 1, \dots, n$; n is number of possible outcomes (or bins)
- Measure of **information** (or **surprise**) $h(X = x_i)$: x_i very probable, i.e., $p(X = x_i)$ is high then $h(X = x_i)$ should be low, and v.v.
- Realizations x_i and y_j of independent X and Y , i.e., the two-dimensional pdf $p(X = x_i, Y = y_j)$ equals the product of the one-dimensional marginal pdfs $p(X = x_i)p(Y = y_j)$, we would like the joint information content to equal the sum of the marginal information contents, i.e., $h(X = x_i, Y = y_j) = h(X = x_i) + h(Y = y_j)$
- Feasible with $h(X = x_i) = -\ln(p(X = x_i))$
- The expectation $H(X)$ of the information content is termed the (Shannon) **entropy**: $H(X) = -\sum_{i=1}^n p(X = x_i) \ln(p(X = x_i)) = -E\{\ln(p)\}$

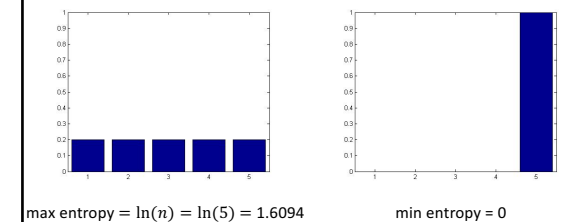
Entropy

- $\lim_{p \rightarrow 0+} p \ln(p) = 0$
- Unit: nat, nit or nepit (natural log); bit (log base 2); ban or dit (log base 10)
- For continuous stochastic variable $H(X) = -\int p(x) \ln(p(x)) dx$
 - Differential entropy; since $p(x)$ can be >1 , $H(X)$ can be negative (or infinite)
 - Gaussian pdf has max entropy for continuous distributions with finite variance, $H(X) = \{1 + \ln(2\pi\sigma^2)\}/2$; $H(X) = \ln[(2\pi e)\mathbb{I}]/2$
- Discrete X : be careful if binning is applied
- Empirical entropy $\hat{H}(X) = -\sum_{i=1}^N \ln(p(X = x_i)) / N$; N is number of obs
- Rényi entropy: $H_\alpha(X) = \frac{1}{1-\alpha} \ln(\sum_{i=1}^n p(X = x_i)^\alpha)$, $\alpha > 0$ and $\alpha \neq 1$
 - $\alpha \rightarrow 1$: Shannon entropy
 - $\alpha = 2$: Collision or **Rényi entropy**: $H_2(X) = -\ln(\sum_{i=1}^n p(X = x_i)^2) = -\ln(E\{p\})$

Entropy

- Shannon entropy**: $-E\{\ln(p)\}$
- Rényi entropy**: $-\ln(E\{p^\alpha\})$ $E\{p\}$ is energy
- Entropy is the average amount of information (or surprise) obtained from obs**
- Entropy is a measure of order**
 - Physical system: heat source/sink, ability to do work, low entropy**
 - Source/sink same temperature: no work done, max entropy**
 - Histogram of temperature of source/sink**
 - A bit like my (low entropy) office: a lot of potential for work**

Entropy



Entropy

- Claude Elwood Shannon (1916-2001), A Mathematical Theory of Communication. *Bell System Technical Journal*, 1948: **application areas communication, compression** – “father of information theory”
- Alfréd Rényi (1921-1970), contributions in combinatorial, graph theory, number theory but mostly in probability theory
- John von Neumann (1903-1957), major contributions in mathematics, physics, economics, computer science, and statistics (source http://en.wikipedia.org/wiki/John_von_Neumann):
use the term entropy not only because of the link to physics, but also because “nobody knows what entropy really is, so in any discussion you will always have an advantage”, (Christopher Michael Bishop, 2006)

Entropy

- Early use in maximum entropy (image) reconstruction (B. Roy Frieden, 1972)
 - Entropy as regularizer in inverse problems to
 - suppress low-intensity ripple
 - sharpen point sources
 - in e.g. astronomical data
- Equality of distributions of salaries/income
 - Gini index – first order series expansion of $\ln(x)$ in Shannon entropy around $x = 1$, $\ln(x) \cong x - 1$: $\text{Gini} = \sum_{i=1}^n p(X = x_i)(1 - p(X = x_i)) = E\{1 - p\}$
- None of this described further here, on to simpler things

Entropy of Posterior Probability

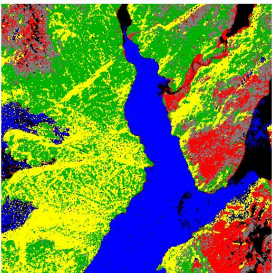


Igaliko
Brattahlíð (Qassiarsuk)
Landsat MSS

Entropy of Posterior Probability

ML, quadratic classification, Mahalanobis reject

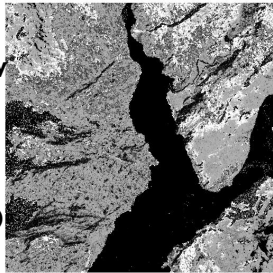
ML classification with reject class



Igaliko
Brattahlíð (Qassarsuk)
Landsat MSS

Entropy of Posterior Probability

Entropy of posterior probability to characterize certainty of classification




Igaliko
Brattahlíð (Qassarsuk)
Landsat MSS

Spectral Information Measure

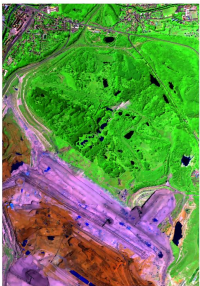
Entropy of DNs to characterize spectral information content

Chen-I. Chang, 2000



Igaliko
Brattahlíð (Qassarsuk)
Landsat MSS


Spectral Information Measure



Sokolov
HyMap, 110 bands

Spectral Information Measure

Entropy of DNs to characterize spectral information content

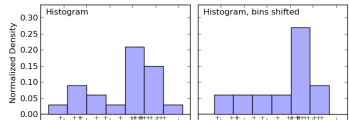


Sokolov
HyMap, 110 bands

Density Estimation

- Discrete X : be careful if binning is applied
 - The estimated histogram is not smooth and it depends on the end points of bins and the width of bins
 - Using kernel density estimators where we center a kernel on each observation, we may obtain smoother histograms that do not depend on bin end points
- Kernel density estimator value t is $\sum_{i=1}^N \varphi(t - x_i)$

$$\varphi(t) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{t^2}{2\sigma^2}\right)$$
 - Issue: determine bandwidth
- This parameterization



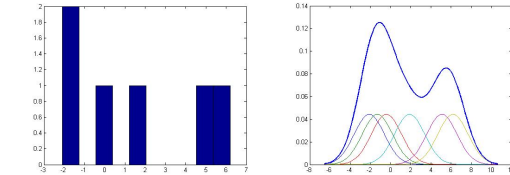
Density Estimation

- Discrete X : be careful if binning is applied
 - The estimated histogram is not smooth and it depends on the end points of bins and the width of bins
- Using kernel density estimators where we center a kernel on each observation, we may obtain smoother histograms that do not depend on bin end points
- Kernel density estimator (Parzen window estimator) for the pdf of X at value t is $\sum_{i=1}^N \varphi(t - x_i)/N$; often with Gaussian kernel

$$\varphi(t) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{t^2}{2\sigma^2}\right) \quad (\text{Bernard W. Silverman, 1986})$$
 - Issue: determine bandwidth σ
- This parameterization facilitates optimization

Density Estimation

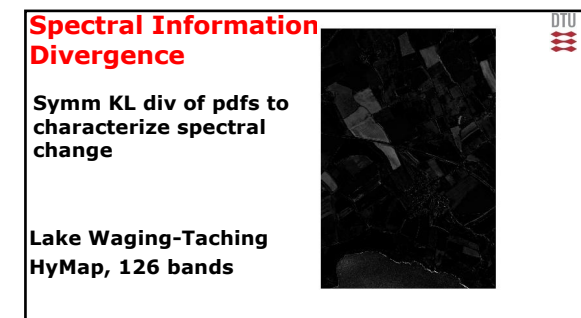
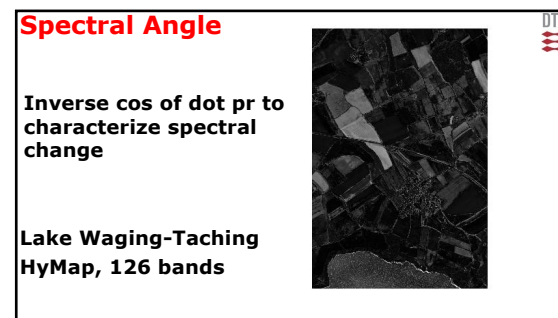
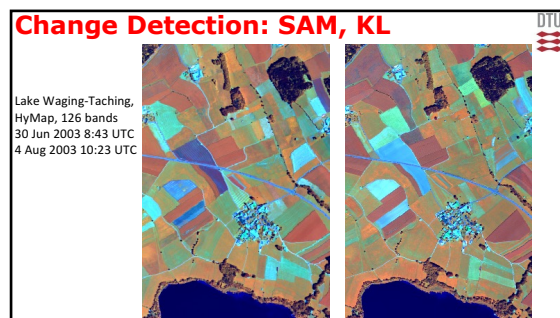
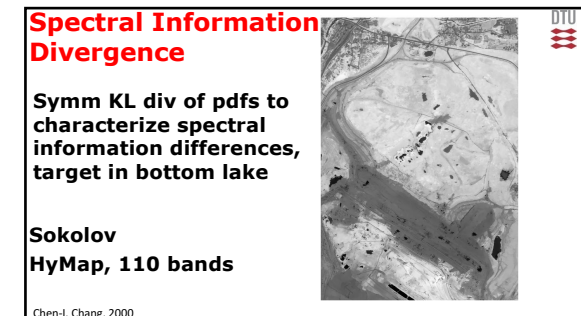
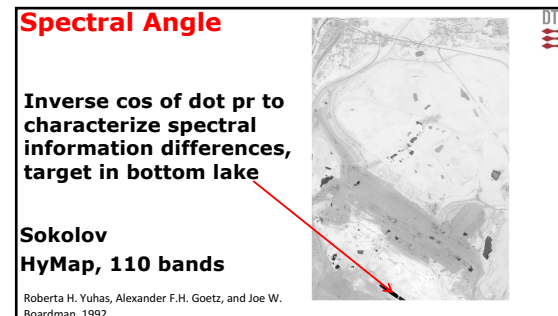
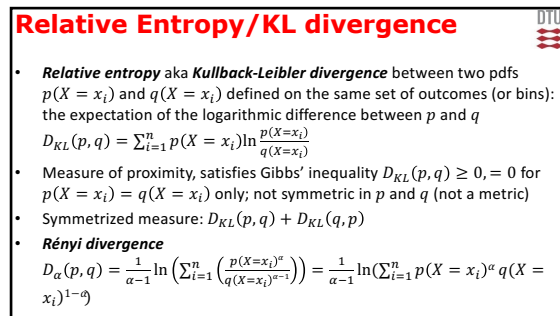
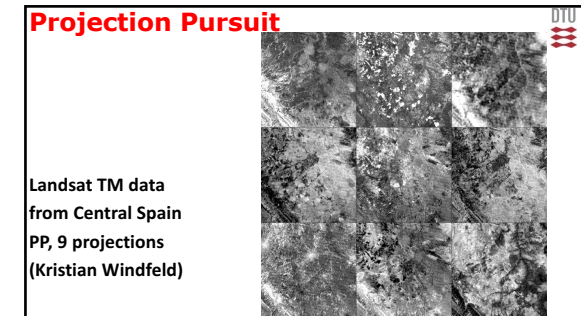
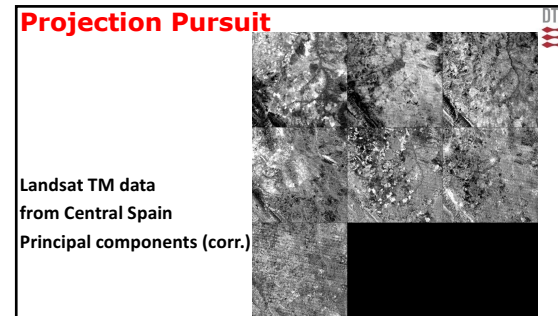
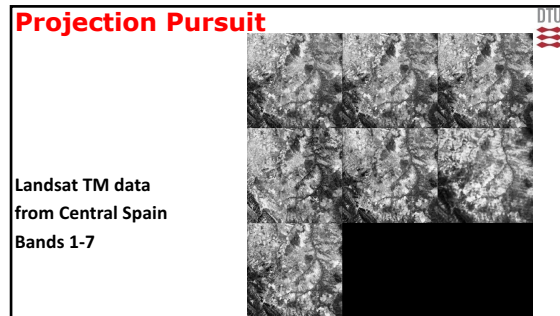
$x = [-2.1 \ -1.3 \ -0.4 \ 1.9 \ 5.1 \ 6.2]$



Source: http://en.wikipedia.org/wiki/Kernel_density_estimation

Projection Pursuit

- PCA (Harold Hotelling, 1933): linear combinations $a^T X$ of m (de-meanned) variables in X with max variance
 - Higher order components orthogonal to lower order components
- PP (Jerome H. Friedman and John Tukey, 1974): linear combinations $a^T X$ of m (de-meanned) variables in X with max deviation from uniform/Gaussian distribution after mapping with essentially Φ (standard normal cdf) measured by integral-squared distance or with min entropy
 - Higher order components: replace structure in projections with Gaussian noise and go again, i.e., no orthogonality, more than m projections
- Quite similar to ICA (Pierre Comon, 1994)



Mutual Information

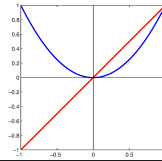
- **Kullback-Leibler divergence**

$$D_{KL}(p, q) = \sum_{i=1}^n p(X = x_i) \ln \frac{p(X = x_i)}{q(X = x_i)}$$

- **Mutual information:** measure of extent to which X and Y are independent, KL divergence between $p(X = x_i, Y = y_j)$ and $p(X = x_i)p(Y = y_j)$

$$I(X, Y) = \sum_{i,j} p(X = x_i, Y = y_j) \ln \frac{p(X = x_i, Y = y_j)}{p(X = x_i)p(Y = y_j)}$$

- $I(X, Y)$ symmetric in X and Y (a metric)
- $I(X, Y) = H(X) + H(Y) - H(X, Y)$
- A normed version as supplement to correlation

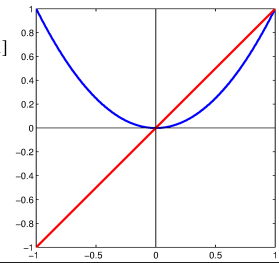


Mutual Information

- CCA: linear combinations $a^T X$ and $b^T Y$ of m and k (de-meant) variables in X and Y with max correlation (Harold Hotelling, 1936)
 - Higher order components orthogonal to lower order components
- **Replace correlation in CCA with mutual information** (Xiangrong Yin, 2004; Karasuyama & Sugiyama, 2012; Vestergaard & Nielsen: CIA, 2015)
 - Higher order components: min MI with lower order components; no orthogonality; more than $\max(m, k)$ components
 - CIA handles large datasets
- [Examples](#)
- [CIA paper](#)

Toy example

- $X = [x_1 \ x_2]$, $Y = [y_1 \ y_2]$
- x_1 sampled equidistantly on $[-1 \ 1]$
- $y_1 = x_1^2 + \varepsilon_1$
- $x_2 = \varepsilon_2$
- $y_2 = \varepsilon_3$
- $\varepsilon_i \sim N(0, 0.1^2)$
- $\max I(a^T X, b^T Y)$



Toy example

