



“A breakthrough in Machine Learning
would be worth ten Microsofts.”
– Bill Gates (2004)

Lectures 1&2: Introduction and Course Overview

Prof. Krishna R. Pattipati
BoT Distinguished Professor &
UTC Chair Professor in Systems Engineering
Dept. of Electrical and Computer Engineering
University of Connecticut

Contact: krishna.pattipati@uconn.edu; (860) 486-2890

Fall 2018

August 27 & September 10, 2018



Reading List

- ❑ Chapters 1 & 2 of Bishop
- ❑ Chapter 1 of *Model-based Machine Learning* by Winn and Bishop (Online)
- ❑ Chapters 1 & 2 of Murphy
- ❑ Chapter 2 of Theodoridis
- ❑ Class Notes
- ❑ Science:
 - Special Issue on AI, July 17, 2015, Pages 248-278
 - AI Transforms Science, July 7, 2017, Pages 16-30
- ❑ Web
 - <https://waitbutwhy.com/2015/01/artificial-intelligence-revolution-1.html>
 - <https://waitbutwhy.com/2015/01/artificial-intelligence-revolution-2.html>
- ❑ YouTube (Bishop's talk on AI: The History and Future)
 - https://www.youtube.com/watch?v=8FHBh_OmdsM



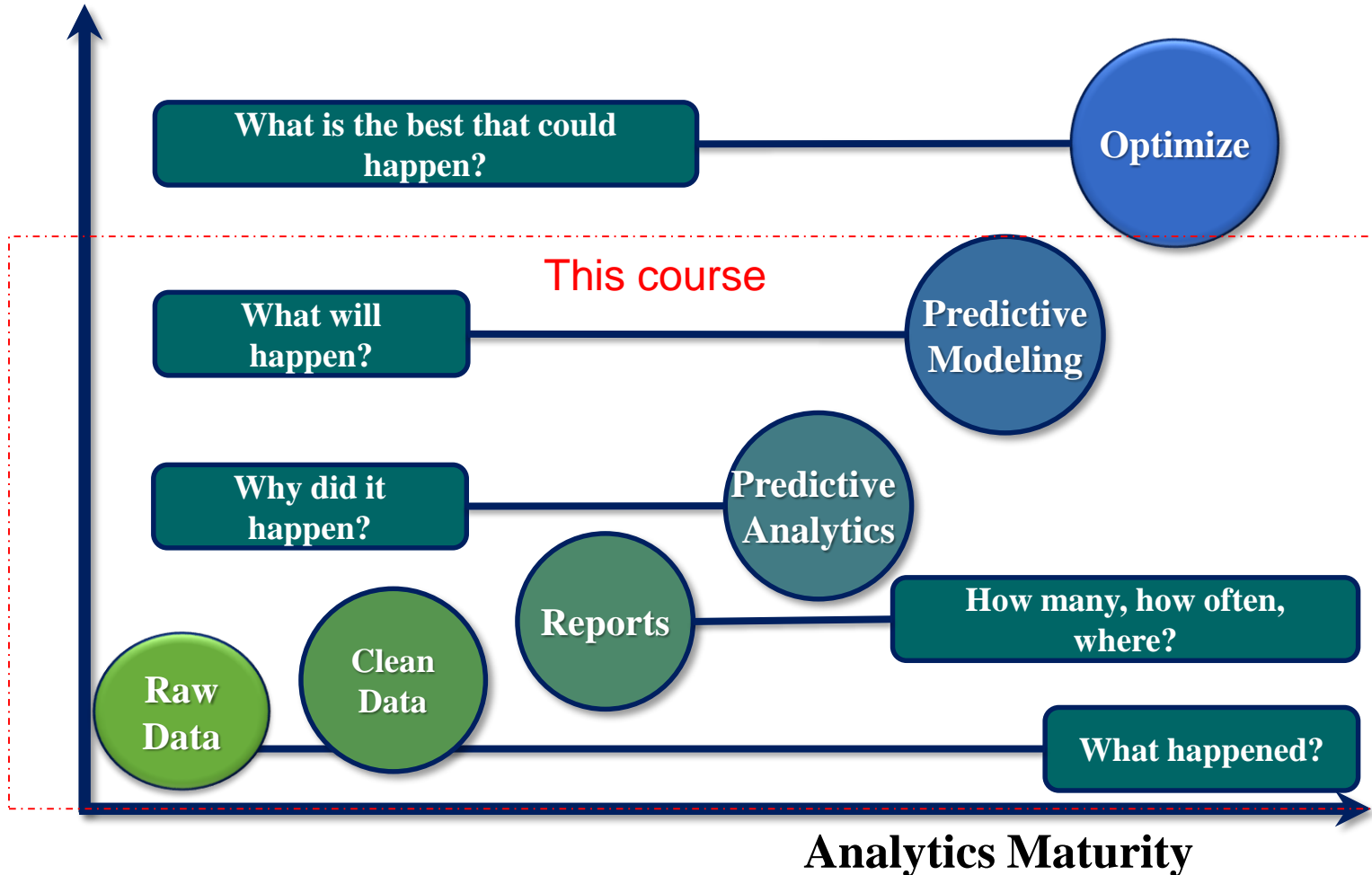
Why Do Machines Need to Learn?

- ❑ Deluge of data (trillions of web pages, 10 years of video content in a single day on YouTube, retail sales, Facebook, twitter, blogs,...)
- ❑ Don't know Concise I/O or cause-effect relationships
- ❑ Are there hidden relationships and correlations (“Patterns”) in large data?
- ❑ Need to adapt to an unknown environment (“Uncertainty”)
- ❑ New knowledge is constantly being rediscovered (“Dynamic”)
- ❑ Entering knowledge by humans is tedious



Why ML? : Data to Decisions

Operational Advantage





Machine Learning Applications

- ❑ Computational Biology
 - Classifying protein sequences, identifying protein-coding genes, genetic networks from gene expression data, ...
- ❑ Medical Imaging
 - Computer-aided diagnosis, image-guided therapy
 - EEG classification, multi-model image fusion
- ❑ Information Retrieval/Natural Language Processing
 - Text/audio/image retrieval
 - Parsing/translation/text analysis/document classification/spam filters
- ❑ Robotics
 - Autonomous driving, planning, control
- ❑ Speech processing
 - Voice identification, speech recognition, speech → text → speech (in another language)
- ❑ Recommender Systems
 - Collaborative filtering, market basket analysis,....
- ❑ Scientific Data Analysis
 - E.g., NASA missions, Hadron Collider,...
- ❑ Financial Prediction and Automated Trading
- ❑ Computer Games, Gambling,....



What is Learning?

- Learning: *Improving with Experience at some Task*
 - *Improve over task, T*
 - *With Respect to Performance Measure, P*
 - *Based on Experience, E*

- Example: *Spam Filtering* (supervised learning)
 - *Spam: An E-mail that the user does not want to receive and has not asked to receive*
 - *T : Identify Spam E-mails*
 - *P : % Spam E-mails that were filtered (Detection, Sensitivity)*
 - % of ham (non-spam) E-mails that were incorrectly filtered-out (False Alarms)*
 - *E : A Database of E-mails that were labelled by users/experts*

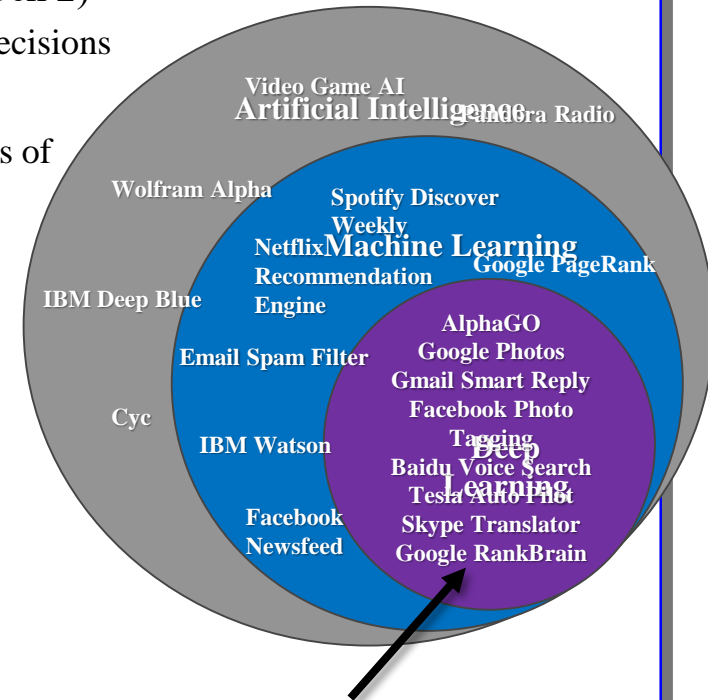
- Learning spans multiple fields: *EE, CS, Statistics, Cognitive Science, Computational Neuroscience, Computational Biology, Social Sciences, Financial Engineering, Economics,*



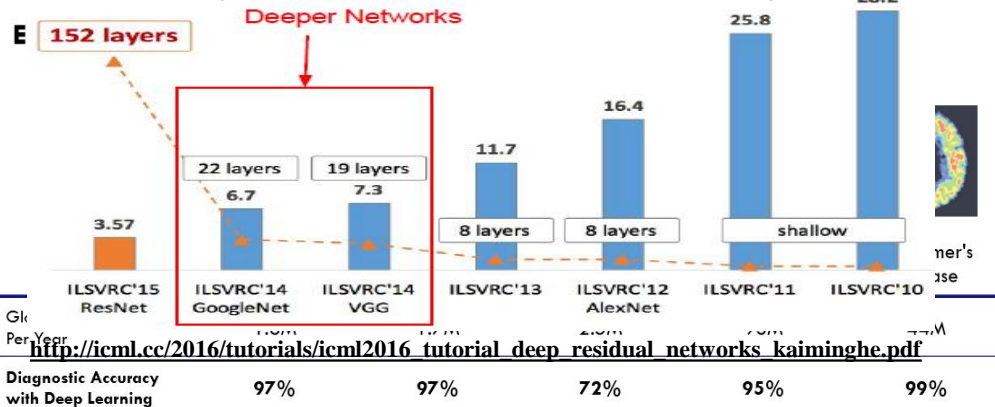
AI versus ML versus Deep Learning

AI: The study of the Computations that make it possible to **Perceive, Reason and Act**

- Classic AI is based on deductive logic (Gen 1)
 - Rules are based on human ingenuity
- Machine Learning is based on learning “models” from data (Gen 2)
 - Rules/Models derived from data are used for predictions and decisions
- Deep learning seeks to “mimic” the biological brain (Gen 3)
 - Learn and recognize patterns of input data using multiple layers of abstraction
 - Applications:
 - Image classification: Facial recognition
 - Speech recognition in phones: Siri, Google assistant
 - Game playing: Chess, Atari, Alpha GO
 - Disease prediction
 - Self-driving cars



ImageNet Classification Top-5 Error in Percentage



“Greedy, Brittle, Opaque and Shallow” Deep learning addresses the so-called ‘1-second’ problem (narrow intelligence)



What Made Deep Learning Feasible?

- ❑ Deluge of data and ability to crowd source the labeling process for supervisory learning ...Availability of *large scale (and often pre-trained) labeled images* (e.g., Alexnet, ImageNet, VGG16) and advances in *cross-modal* (image, text/data, video and audio) *representations*.... “*Big Data*”
- ❑ Availability of *Graphics Processing Units* (GPUs) that enable training of very large image/text/speech datasets from several weeks to a few days
- ❑ *New nonlinearities* (e.g., rectified linear units) and *signal normalization* that avoid numerical problems associated with gradient computations
- ❑ *Convolution and max-pooling* operations that exploit local connectivity to reduce the dimension of the weight space, control overfitting and make the network robust
- ❑ Concept of *dropout* to realize an exponentially large ensemble of networks from a single network
- ❑ Advances in *stochastic optimization, adversarial training and regularization* for robust network training
- ❑ *More layers capture more invariances* (Information-theoretic insights)
- ❑ *Much better software tools* (Keras, TensorFlow, Theano, Torch,...)
- ❑ **Result:** *Features learned rather than hand-crafted, better accuracy (but easily fooled!)*



Types of Learning

□ Learning from Data: Model Parameters, W (weights)

$$W = f(D) \quad D \rightarrow \text{data used to learn (training data)}$$

There are at least four forms of learning (others: active, transfer)

- Supervised Learning (or learning with a teacher)
- Unsupervised Learning
- Semi-supervised Learning
- Reinforcement Learning (like learning with a critic)

a) Supervised Learning:

$$w = f[\{\underline{x}_n, \underline{z}_n\}_{n=1}^N] \quad (\text{know inputs and correct outputs})$$

- Classification
- Regression
- Time series prediction

b) Unsupervised Learning:

$$w = f[\{\underline{x}_n\}_{n=1}^N] \quad (\text{has input data only})$$

- Clustering; Density Estimation
- Outlier Detection
- Compression

c) Semi-supervised Learning: Practical setting; Active learning

$$w = \{f[\{\underline{x}_n, \underline{z}_n\}_{n=1}^{N_1}]; \{f(\underline{x}_n)\}_{n=N_1+1}^N\} \quad (N_1 \ll N: \text{ Few outputs are labeled, but most are not})$$

d) Reinforcement Learning: Useful in adaptive control and decision making

$$w = f[\{\underline{x}_n\}_{n=1}^N; \text{ correctness/effectiveness of outputs}] \quad (\text{learning with a critic})$$



Things to Remember in ML Model Building

- ❑ Make Sure Training Data is Adequate (“Data Matters More than Algorithms”)
- ❑ Make Sure Data is Representative (“No sampling Bias”)
- ❑ Clean up the Data (“Detect and Fix/Remove Outliers”, “Ignore/Estimate Missing Data”, “Minimize Noise Effects in the Data”)
- ❑ Feature Engineering (“Feature Selection”, Feature Extraction”, “Seek new Information”)
- ❑ Avoid Overfitting and Underfitting (“Kiss” Principle, “Occam’s Razor”, “Theory of Parsimony”, “Bias-Variance Tradeoff”)
- ❑ Understand the difference between Training Error and Generalization (Out-of-sample) Error (“Split Data into Training, Validation and Test Data Sets”, “No data snooping”, “Test and Validate the Model”, “Cross Validation”, “Bootstrap”, “Model Selection”, “Bayesian Model Averaging”)
- ❑ Exploit Domain Knowledge (“Relationships and Constraints”, “No Free Lunch Theorem”)



Some Interesting Datasets

Dataset	Accessible Location
IRIS	http://www.statlab.uni-heidelberg.de/data/iris/
MNIST	http://yann.lecun.com/exdb/mnist/
20Newsgroups	http://cs.nyu.edu/~roweis/data.html
Olivetti Faces	http://cs.nyu.edu/~roweis/data.html
Alarm Network	http://www.bnlearn.com/bnrepository/

UCI Machine Learning Repository: <http://archive.ics.uci.edu/ml/index.php>

Kaggle Data Sets: <https://www.Kaggle.com/datasets>

Amazon AWS: <http://aws.amazon.com/fr/datasets>

Open Source Data Portals: <http://dataportals.org> ; <http://opendatamonitor.eu> ;
<http://quandl.com>;

Wiki: https://en.wikipedia.org/wiki/List_of_datasets_for_machine_learning_research

Deep Learning Datasets: <http://deeplearning.net/datasets/>

Blog: <https://blog.algorithmia.com/machine-learning-datasets-for-data-scientists/>



Fisher-Anderson's Iris Dataset



iris setosa



iris versicolor

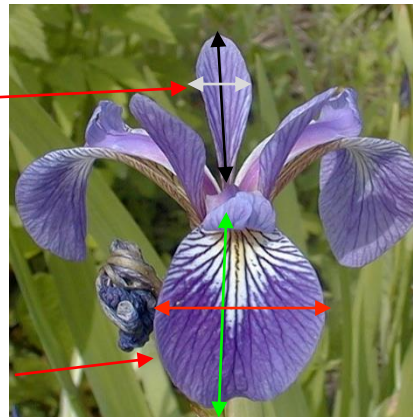


iris virginica

What type of measurements will help distinguish these flowers well?

petal

sepal



- ↔ petal length
- ↔ petal width
- ↔ sepal length
- ↔ sepal width

<http://www.statlab.uni-heidelberg.de/data/iris/>

The answer to “what measurements to take” typically comes from subject matter experts, Botanists in this case.



Sample Data of Iris Dataset

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5	3.6	1.4	0.2	setosa
7	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.9	3.1	4.9	1.5	versicolor
5.5	2.3	4	1.3	versicolor
6.5	2.8	4.6	1.5	versicolor
6.3	3.3	6	2.5	virginica
5.8	2.7	5.1	1.9	virginica
7.1	3	5.9	2.1	virginica
6.3	2.9	5.6	1.8	virginica
6.5	3	5.8	2.2	virginica

Both petal length and width clearly distinguish setosa; but the other two are harder to distinguish

1. R. A. Fisher (1936). "The use of multiple measurements in taxonomic problems." *Annals of Eugenics*. 7 (2): 179–188.
2. Edgar Anderson (1936). "The species problem in Iris." *Annals of the Missouri Botanical Garden*. 23 (3): 457–509.



Get to Know the Data

- ❑ Continuous Features: For each feature, tabulate
 - Count
 - % Missing
 - Cardinality
 - Minimum and Maximum
 - 1st and 3rd Quartiles
 - Mean and Median
 - Variance and Standard Deviation
 - Skewness
 - Kurtosis
- ❑ Discrete (Categorical) Features
 - Count
 - % Missing
 - Cardinality
 - First and Second Mode, Mode Frequency and % Mode Frequency



Data Statistics: Mean

□ Data Matrix X : an N by p matrix

- N = Number of data points (data samples)
- p = Number of features (e.g., attributes, measurements, variables)

□ Central Tendencies

- Mean of each column feature (MATLAB: `mean(X)`)

$$\underline{\hat{\mu}} = \frac{X^T \underline{e}}{N}; \quad \underline{e} \sim N \text{ column vector of } 1^s$$

- Iris Dataset:

- Overall mean: $\underline{\mu}^T = [5.8433 \quad 3.0573 \quad 3.7580 \quad 1.1993]$
- Setosa: $\mu_1^T = [5.0060 \quad 3.4280 \quad 1.4620 \quad 0.246]$
- Versicolor: $\mu_2^T = [5.9360 \quad 2.7700 \quad 4.2600 \quad 1.3260]$
- Virginica: $\mu_3^T = [6.5880 \quad 2.9740 \quad 5.5520 \quad 2.0260]$

- Setosa has smaller petal length and petal width. Easy to classify
- Versicolor has smaller petal width than Virginica; somewhat harder to distinguish between Versicolor and Virginica

Mean Removed Data Matrix

$$X_1 = X - \underline{e} \underline{\hat{\mu}}^T = \underbrace{\left(I - \frac{\underline{e} \underline{e}^T}{N}\right)}_{\text{Projection Matrix, } P} X$$

Note: $P^n = P; n = 2, 3, \dots$



Data Statistics: Median

□ Central Tendencies

➤ Median (median (X))

- Sort the values in each column j in increasing order: $x_{[1],j} \leq x_{[2],j} \leq \dots \leq x_{[N],j}$
- If N is odd

$$Md_j = x_{\lfloor \frac{N+1}{2} \rfloor, j}$$

- If N is even

$$Md_j = \frac{1}{2} [x_{\lfloor \frac{N}{2} \rfloor, j} + x_{\lfloor \frac{N}{2} \rfloor + 1, j}]$$

➤ Iris Dataset

- Overall Median: [5.8000 3.0000 4.3500 1.3000]
- Setosa: [5.0000 3.4000 1.5000 0.2000]
- Versicolor: [5.9000 2.8000 4.3500 1.3000]
- Virginica: [6.5000 3.0000 5.5500 2.0000]

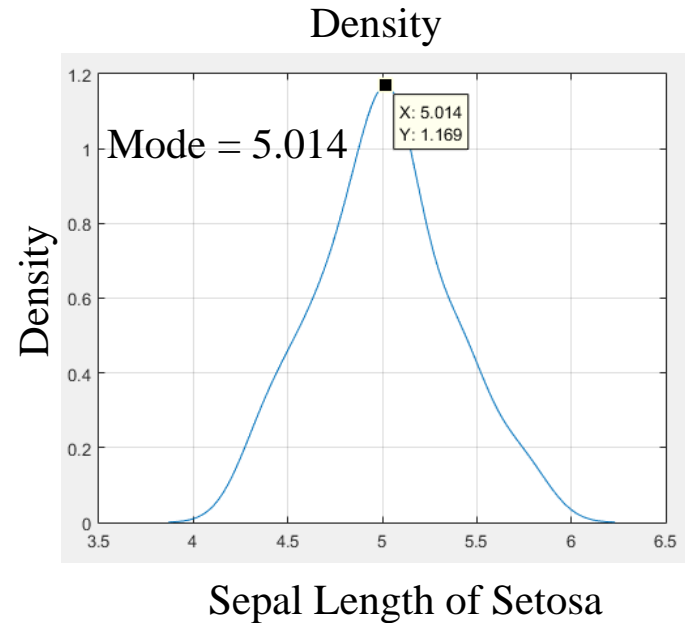
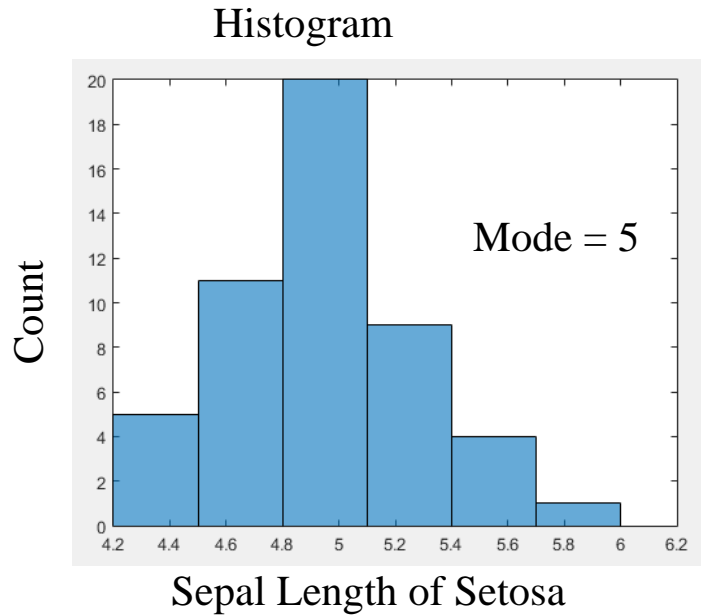
➤ Similar conclusions as with mean



Data Statistics: Mode

Central Tendencies

- Mode of each feature \equiv most frequent value in each column
 - Good for discrete features: $(\text{mode}(X))$
 - For continuous features, do a histogram or density estimate first and calculate the *peak* of the estimate ($\text{histogram}(X(:,j))$) or $\text{ksdensity}(X(:,j))$; $j=1,2,\dots,p$)





Data Statistics: Variations -1

Variations

- Minimum: ($\min(X)$)

$$x_{j,\min} = \min_{1 \leq i \leq N} (x_{ij}); \quad j = 1, 2, \dots, p$$

- Maximum: ($\max(X)$)

$$x_{j,\max} = \max_{1 \leq i \leq N} (x_{ij}); \quad j = 1, 2, \dots, p$$

- Range: ($\text{range}(X)$)

$$r_j = x_{j,\max} - x_{j,\min}; \quad j = 1, 2, \dots, p$$

- Sample Variance for *iid* samples: ($\text{var}(X)$)

$$\hat{\sigma}_j^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \hat{\mu}_j)^2; \quad \hat{\mu}_j = \frac{1}{N} \sum_{i=1}^N x_{ij}; \quad j = 1, 2, \dots, p$$

Why?

You can also compute class-conditional variance

- An unbiased estimator of the covariance matrix (class-conditioned)

Computed from data of each class

$$\hat{\Sigma} = \frac{X_{1c}^T X_{1c}}{n_c - 1} = \frac{X_c^T (I - \frac{e e^T}{n_c}) X_c}{n_c - 1}$$

For common (class-independent) covariance matrix across C classes, replace $(N-1)$ by $(N-C)$ for unbiasedness! See lectures 4-6.

iid = independent and identically distributed

Bernoulli: $\sigma_j^2 = p(1-p)$
 $P(x=1) = p$

Show: $E[\hat{\sigma}_j^2] = \sigma_j^2$
 $\sigma_j^2 = \text{var}(x_{ij})$



Data Statistics: Variations - 2

Variations

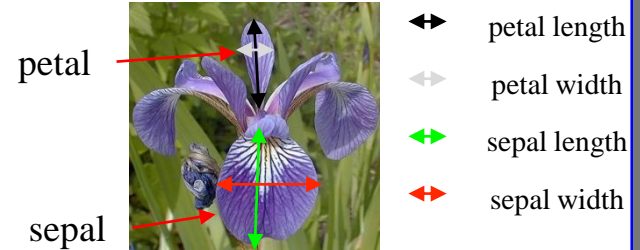
Correlation Coefficient Matrix ($corrcoef(X)$)

$$C_r = [Diag(\hat{\Sigma})]^{-1/2} \hat{\Sigma} [Diag(\hat{\Sigma})]^{-1/2}$$

Iris Dataset

$$\hat{\Sigma}_{versicolor} = \begin{bmatrix} 0.2664 & 0.0852 & 0.1829 & 0.0558 \\ 0.0852 & 0.0985 & 0.0827 & 0.0412 \\ 0.1829 & 0.0827 & 0.2208 & 0.0731 \\ 0.0558 & 0.0412 & 0.0731 & 0.0391 \end{bmatrix}$$

$$C_{r,versicolor} = \begin{bmatrix} 1.0000 & 0.5259 & 0.7540 & 0.5465 \\ 0.5259 & 1.0000 & 0.5605 & 0.6640 \\ 0.7540 & 0.5605 & 1.0000 & 0.7867 \\ 0.5465 & 0.6640 & 0.7867 & 1.0000 \end{bmatrix}$$



<http://www.statlab.uni-heidelberg.de/data/iris/>

1st and 3rd Quartiles of feature j

- Sort data in increasing order: $x_{[1],j} \leq x_{[2],j} \leq \dots \leq x_{[N],j}$
- The *lower half* of data: $x_{[1],j} \leq x_{[2],j} \leq \dots \leq x_{[\frac{N-1}{2}],j}$ for N odd; $x_{[1],j} \leq x_{[2],j} \leq \dots \leq x_{[\frac{N}{2}],j}$ for N even
- The *upper half* of data: $x_{[\frac{N+1}{2}],j} \leq x_{[\frac{N+1}{2}+1],j} \leq \dots \leq x_{[N],j}$ for N odd; $x_{[\frac{N}{2}+1],j} \leq x_{[\frac{N}{2}+2],j} \leq \dots \leq x_{[N],j}$ for N even
- The *first quartile* $Q_{1,j}$ is the median of the *lower half* of the data. This means that about 25% of the numbers in the data set lie below $Q_{1,j}$ and about 75% lie above $Q_{1,j}$
- The *third quartile*, denoted $Q_{3,j}$, is the median of the *upper half* of the data set. This means that about 75% of the numbers in the data set lie below $Q_{3,j}$ and about 25% lie above $Q_{3,j}$.

$$\text{For } N=9: Q_{1,j} = \frac{x_{[2],j} + x_{[3],j}}{2}; Q_{3,j} = \frac{x_{[7],j} + x_{[8],j}}{2}$$

$$\text{For } N=10: Q_{1,j} = x_{[3],j}; Q_{3,j} = x_{[8],j}$$



Data Statistics: Skewness and Kurtosis

□ Bias Corrected Skewness ($skewness(X,0)$)

$$s_j = \left[\frac{\sqrt{N(N-1)}}{(N-2)} \right] \left[\frac{\frac{1}{N} \sum_{i=1}^N (x_{ij} - \hat{\mu}_j)^3}{\left(\frac{1}{N} \sum_{i=1}^N (x_{ij} - \hat{\mu}_j)^2 \right)^{3/2}} \right]; \hat{\mu}_j = \frac{1}{N} \sum_{i=1}^N x_{ij}; j = 1, 2, \dots, p; N \geq 3$$

$$\text{Bernoulli: } s_j = \frac{1-2p}{\sqrt{p(1-p)}} \\ P(x=1) = p$$

- $s_j < 0 \Rightarrow$ left tail is longer; $s_j > 0 \Rightarrow$ right tail is longer; $s_j = 0 \Rightarrow$ symmetric
- Skewness for versicolor: [0.1022 -0.3519 -0.5882 -0.0302]
- Returns from financial stocks have larger-than-Normal skew and kurtosis

□ Bias Corrected Kurtosis ($kurtosis(X,0)$) ... useful for detecting outliers

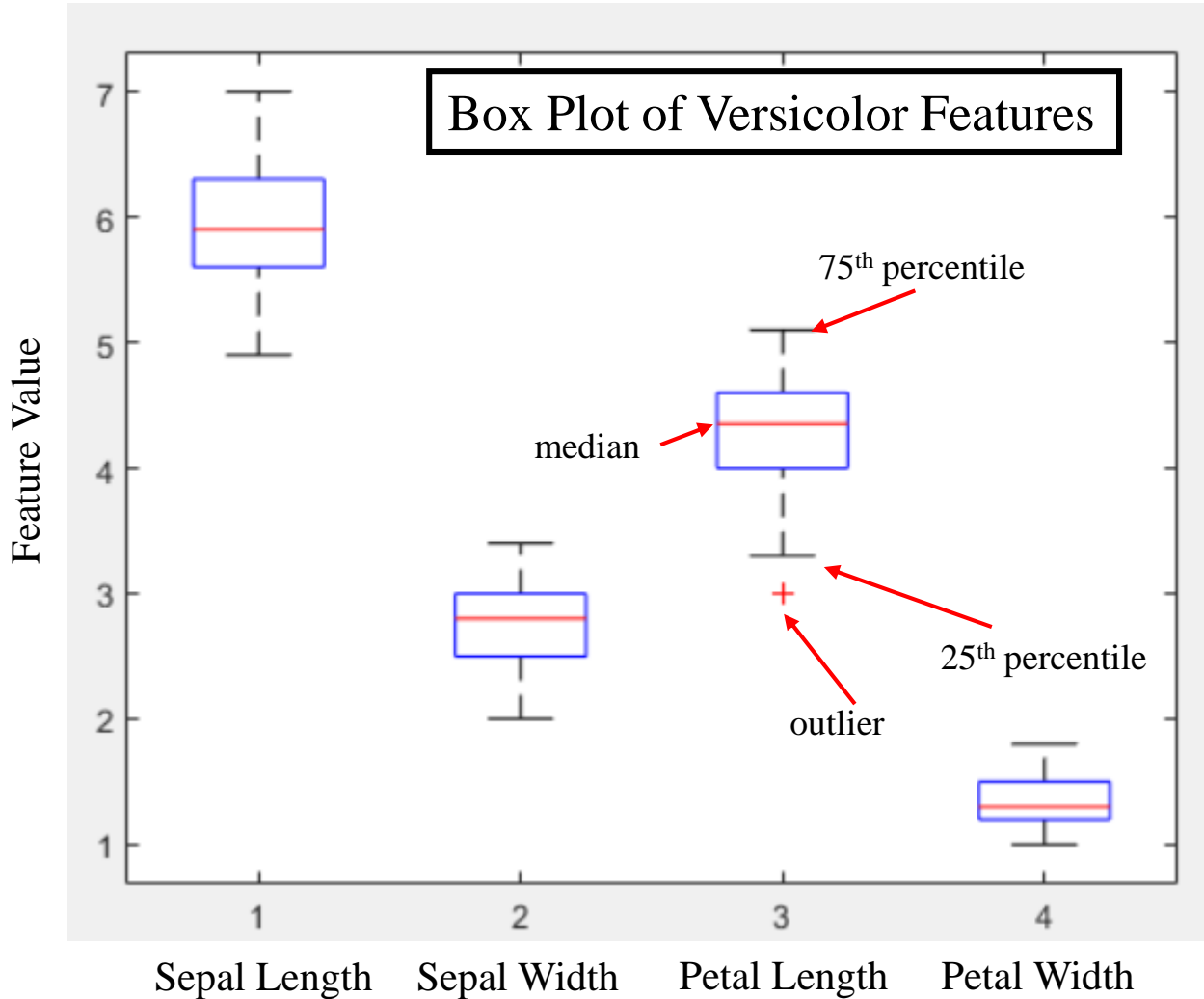
$$\kappa_j = \left[\frac{(N-1)}{(N-2)(N-3)} \right] \left[(N+1) \left(\frac{\frac{1}{N} \sum_{i=1}^N (x_{ij} - \hat{\mu}_j)^4}{\left(\frac{1}{N} \sum_{i=1}^N (x_{ij} - \hat{\mu}_j)^2 \right)^2} \right) - 3(N-1) \right] + 3; \hat{\mu}_j = \frac{1}{N} \sum_{i=1}^N x_{ij}; j = 1, 2, \dots, p; N \geq 4$$

$$\text{Bernoulli: } \kappa_j = \frac{3p^2 - 3p + 1}{p(1-p)}$$

- $\kappa_j = 3 \Rightarrow$ Gaussian and *Mesokurtic* distributions; Excess Kurtosis = $\kappa_j - 3$
- $\kappa_j > 3 \Rightarrow$ *Leptokurtic* or *super-Gaussian* distributions (e.g., Logistic, Laplace, Student-t,...)
- $\kappa_j < 3 \Rightarrow$ *Platykurtic* or *sub-Gaussian* distributions (e.g., Bernoulli with $p=1/2$, uniform)



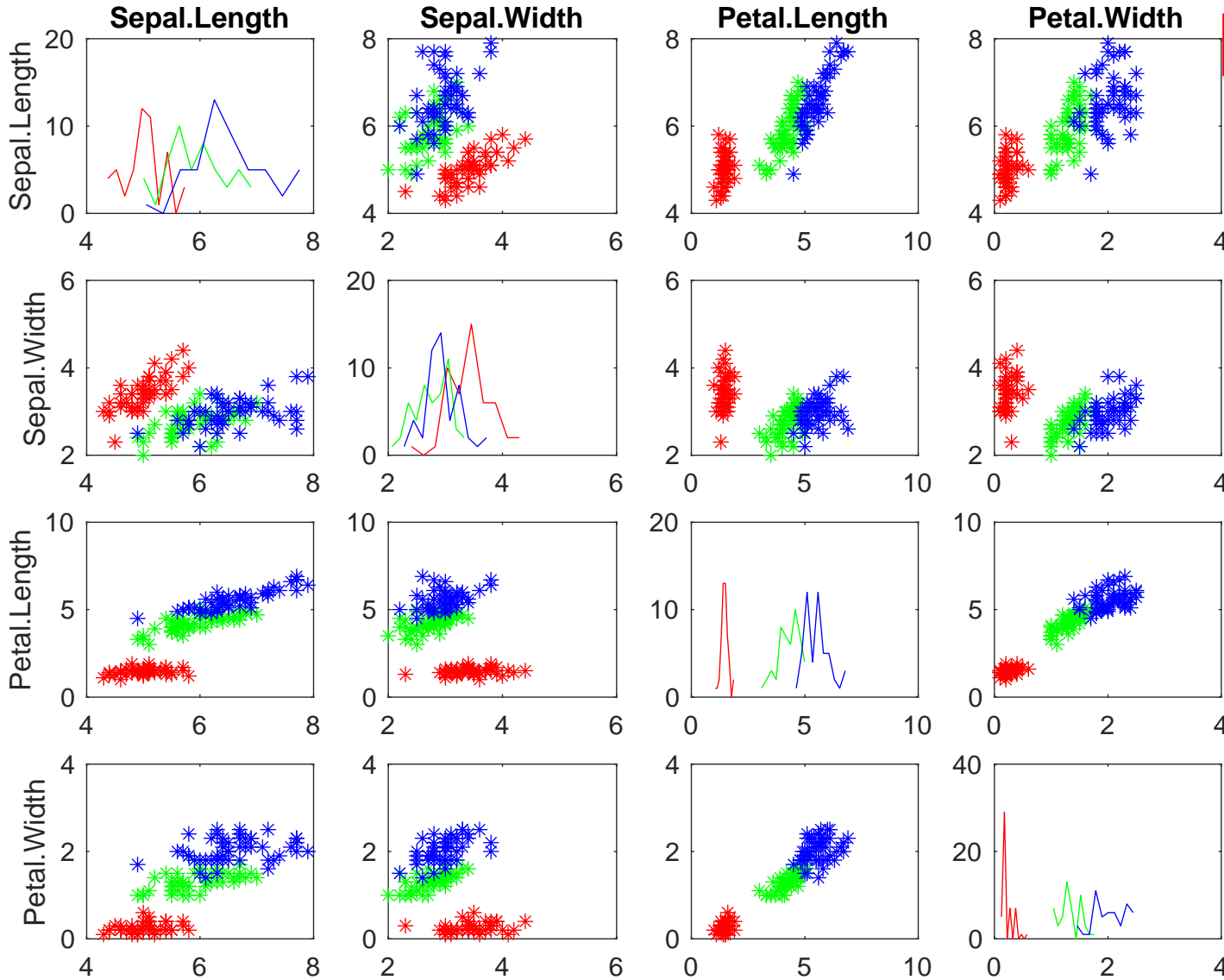
Data Visualization: Box Plot



Bias-corrected Kurtosis: [2.4670 2.6338 3.0479 2.5899]



Very Useful Data Visualization: Matrix Plot



`plotmatrix(X)`

Diagonal:
Histogram/pdf
Off-diagonal:
Scatter or
correlation plot

Colors
Setosa
Versicolor
Virginica

Setosa is
separable



Data Scaling Methods - 1

□ For each Feature

1. *Zero mean and unit variance*

$$x_{ij} \rightarrow s_{ij} = s_{ij}(x_{ij}) = \frac{(x_{ij} - \hat{\mu}_j)}{\hat{\sigma}_j}$$

In Matrix form : $(X - \underline{e}\hat{\underline{\mu}}^T) \text{Diag}(1 / \hat{\sigma}_j)$

2. *Scale it to [0,1] or [-1,1]*

(Often times, $[\varepsilon, 1 - \varepsilon]$ or $[-1 + \varepsilon, 1 - \varepsilon]$) $\varepsilon \approx 0.15 \sim 0.20$ to avoid the saturation of nonlinearities used in neural networks (e.g., sigmoid and tanh))

$$x_{ij} \rightarrow s_{ij} = s_{ij}(x_{ij}) = (1 - 2\varepsilon) \frac{(x_{ij} - x_{\min,j})}{(x_{\max,j} - x_{\min,j})} + \varepsilon \Rightarrow [\varepsilon, 1 - \varepsilon]$$

$$x_{ij} \rightarrow s_{ij} = s_{ij}(x_{ij}) = 2(1 - \varepsilon) \frac{(x_{ij} - x_{\min,j})}{(x_{\max,j} - x_{\min,j})} + (\varepsilon - 1) \Rightarrow [\varepsilon - 1, 1 - \varepsilon]$$



Data Scaling Methods - 2

3. If small values are disproportionately close together and large values are spread out, use *log transformations* to obtain uniformized x^s

$$x_{ij} \rightarrow s_{ij} = s_{ij}(x_{ij}) = \ln \left[\frac{(x_{ij} - x_{\min,j})}{(x_{\max,j} - x_{\min,j})} + 1 \right] / \ln 2$$

4. If large values are close together and small values are further apart, use *exponential transformations* (spreads out large values and pushes small values together)

$$x_{ij} \rightarrow s_{ij} = s_{ij}(x_{ij}) = \frac{1 - e^{-\left[\frac{(x_{ij} - x_{\min,j})}{(x_{\max,j} - x_{\min,j})} \right]}}{1 - e^{-1}}$$

- Coding Classes (Targets): Coding $[1, 2, \dots, C]$ is okay for some (e.g., decision trees)

$$\underline{z}_k = \underline{e}_k = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow z = k \quad \text{or} \quad \underline{z}_k = \begin{bmatrix} -1 / (C - 1) \\ \vdots \\ 1 \\ \vdots \\ -1 / (C - 1) \end{bmatrix} \quad \text{or} \quad \underline{z}_k = \begin{bmatrix} -1 \\ \vdots \\ 1 \\ \vdots \\ -1 \end{bmatrix} \quad \text{Good for MLP}$$



Data Transformation/ Reduction/Cleaning

- Principal Component Analysis (PCA) is the best known and oldest technique for data reduction (Pearson, 1901)
 - Singular Value Decomposition (SVD) if done on data directly
 - Data compression and feature (indicator) selection
- Data Matrix X : an N' by p matrix
 - $N' =$ Number of data points (data samples)
 - $p =$ Number of features (e.g., attributes, measurements)
- Remove *NaN* (Not a Number) elements from the matrix X
 - Use MATLAB function: `X(any(isnan(X), 2), :) = [];`
 - $N =$ Number of data points after removing *NaN* elements
 - New X matrix: an N by p matrix
- Compute the mean of each column feature

$$\underline{\hat{\mu}} = \frac{X^T \underline{e}}{N}; \quad \underline{e} \sim N \text{ column vector of } 1^s$$



Principal Component Analysis Explained

- Remove mean from each column and form a new data matrix X_1

$$X_1 = X - \underline{e}\underline{\hat{\mu}}^T = \left(I - \frac{\underline{e}\underline{e}^T}{N}\right)X$$

- Compute the Covariance Matrix

$$P = \frac{X_1^T X_1}{N-1} = \frac{1}{N-1} \left[X^T X - N \underline{\hat{\mu}}\underline{\hat{\mu}}^T \right] = \frac{1}{N-1} \left[\sum_{n=1}^N (\underline{x}_n \underline{x}_n^T - \underline{\hat{\mu}}\underline{\hat{\mu}}^T) \right]; \text{ a } p \text{ by } p \text{ matrix}$$

- Perform Eigen decomposition of P and select the Eigen vectors corresponding to the top k eigen values \exists

$$\frac{\sum_{i=1}^k \lambda_i}{\text{trace}(P)} > Th \quad (0.95)$$

$$P = V \Lambda V^T;$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \geq \lambda_{k+1} \geq \dots \geq \lambda_p$$

- Reduced Data Matrix

$$X_2 = X_1 V_k \quad \text{an } N \times k \text{ matrix}$$

$$V_k = \text{first } k \text{ columns of } V; \text{ a } p \times k \text{ matrix}$$

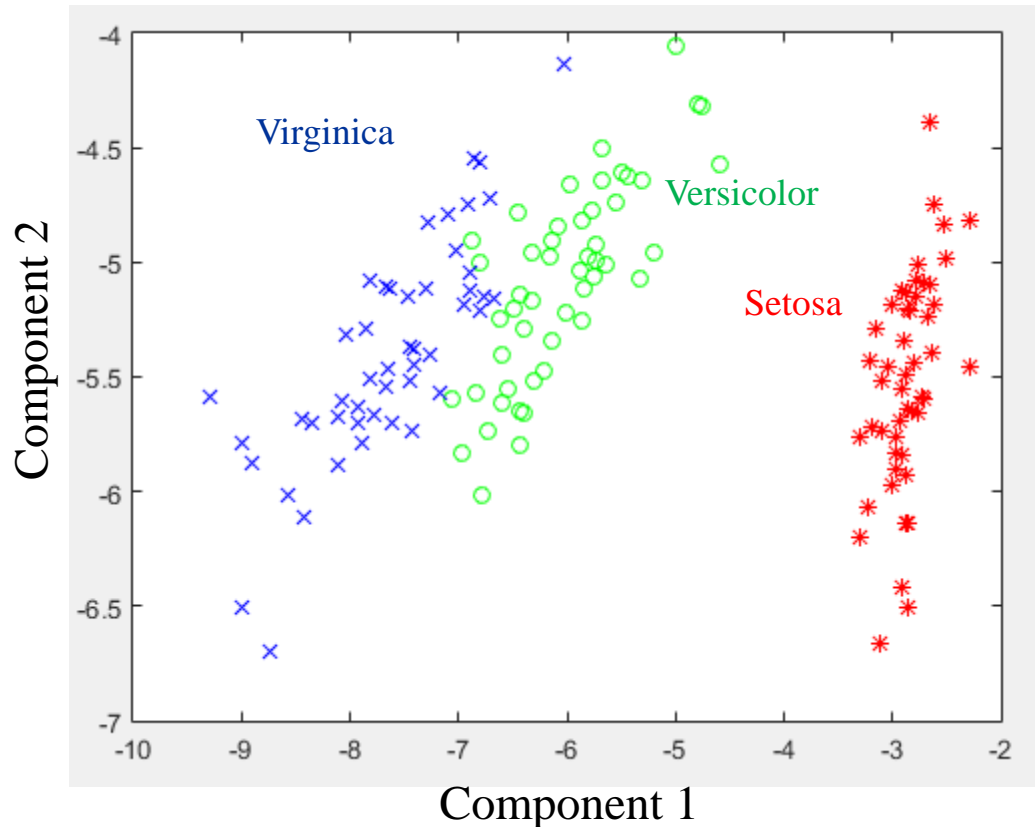
- Residual Matrix and Covariance

$$R = X_1 [I_p - V_k V_k^T] \quad \text{an } N \times p \text{ matrix}; \Sigma_r = [I_p - V_k V_k^T] P [I_p - V_k V_k^T] \text{ a } p \times p \text{ matrix}$$



Iris Data Visualization using PCA

- ❑ First principal component explains 92.46 % of variability in the data
- ❑ First 2 principal components explain 97.8% of variability in the data



- ❑ Setosa is easy to distinguish; some ambiguity between Versicolor & Virginica



Machine Learning Tasks

Machine Learning seeks to solve many generic tasks of interest.

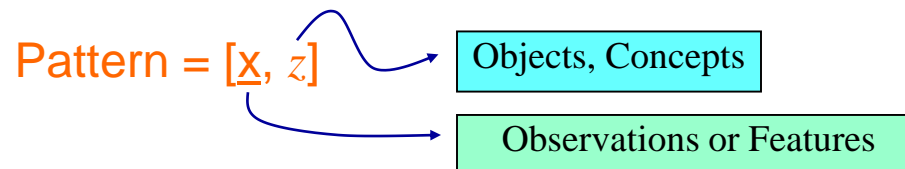
1. Classification
2. Regression / Interpolation
3. Density Estimation
4. Forecasting / Prediction ... basically regression
5. Optimization under uncertainty
6. Adaptive Control
7. Adaptive Communications



Problem 1: Pattern Classification

□ Pattern Classification

- ❖ **Pattern:** Pattern means something exhibiting certain regularities, something able to serve as a model, something representing the concept of what was observed. A pattern is a collection of observations connected in time or space or both.
- ❖ **Classification:** Assign an input pattern (e.g., speech waveform, handwritten character, etc.) represented by a feature vector $\underline{x} = [x_1 \ x_2 \ \dots \ x_p]^T$ to one or more specified classes or concepts



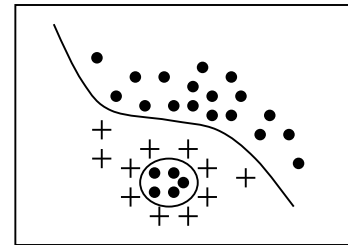
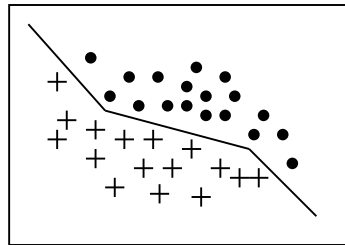
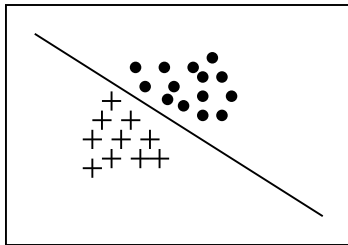
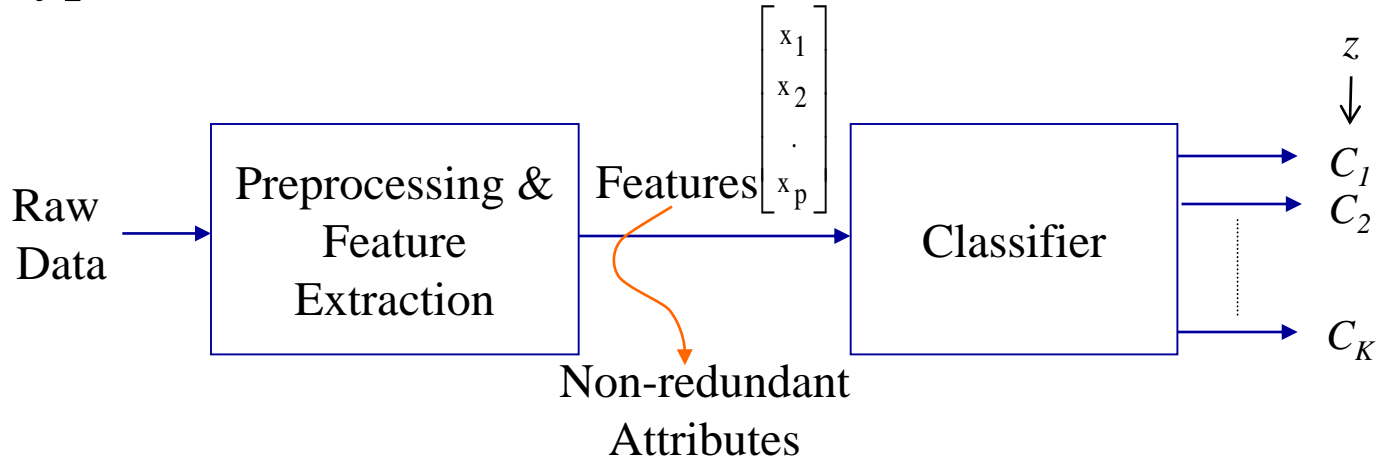
Pattern classification is the process of inferring z from \underline{x} : $\underline{x} \rightarrow z$

\underline{x} can be documents, images, DNA sequences, or graphs,....
 z can be classes, ranks, real values, part of speech tagging, graph,....



Pattern Classification Process

- Typical Process of Pattern Classification





Application Example: Classifying Fish

From Duda, Hart and Stork

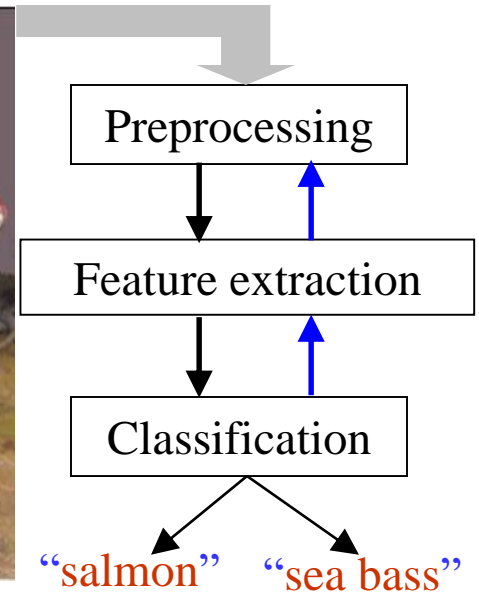


FIGURE 1.1. The objects to be classified are first sensed by a transducer (camera), whose signals are preprocessed. Next the features are extracted and finally the classification is emitted, here either “salmon” or “sea bass.” Although the information flow is often chosen to be from the source to the classifier, some systems employ information flow in which earlier levels of processing can be altered based on the tentative or preliminary response in later levels (blue arrows). Yet others combine two or more stages into a unified step, such as simultaneous segmentation and feature extraction.



Classification using Width and Lightness

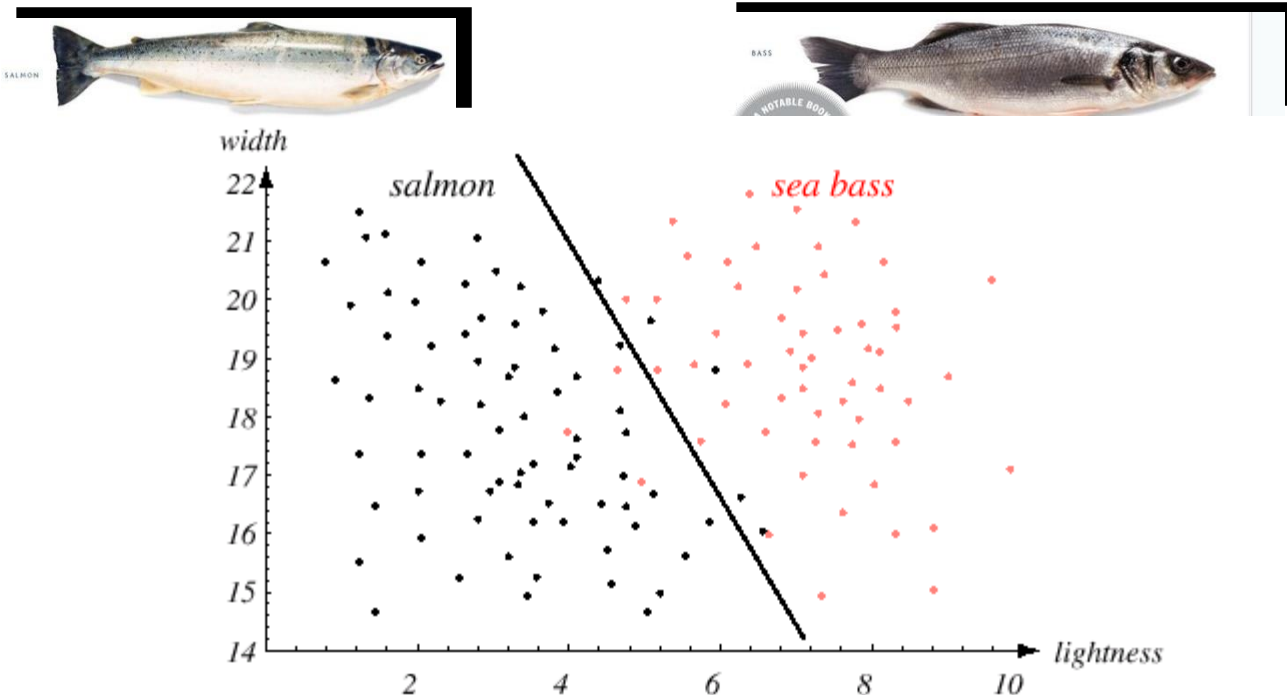


FIGURE 1.4. The two features of lightness and width for sea bass and salmon. The dark line could serve as a decision boundary of our classifier. Overall classification error on the data shown is lower than if we use only one feature as in Fig. 1.3, but there will still be some errors. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Linear (or) Hyperplane Classifier
- Simple, but not necessarily the best for this problem

From Duda, Hart and Stork

How to Select the Right Decision Boundary?

From Duda, Hart and Stork

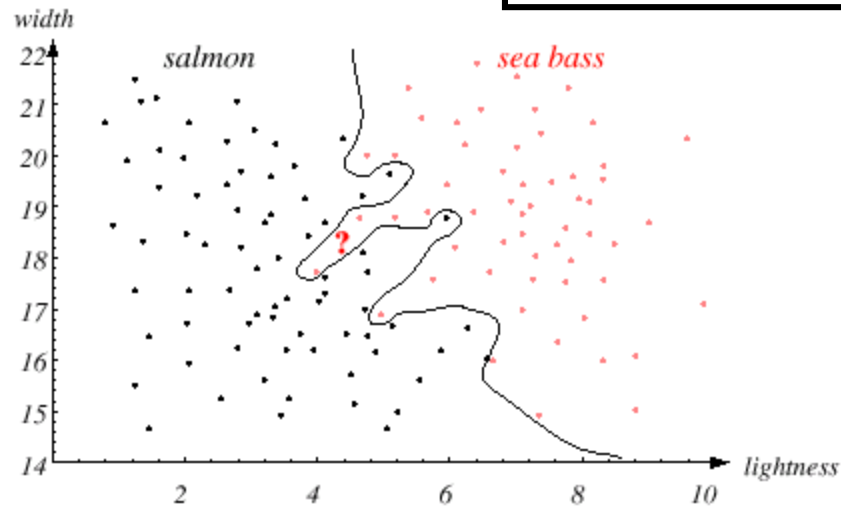
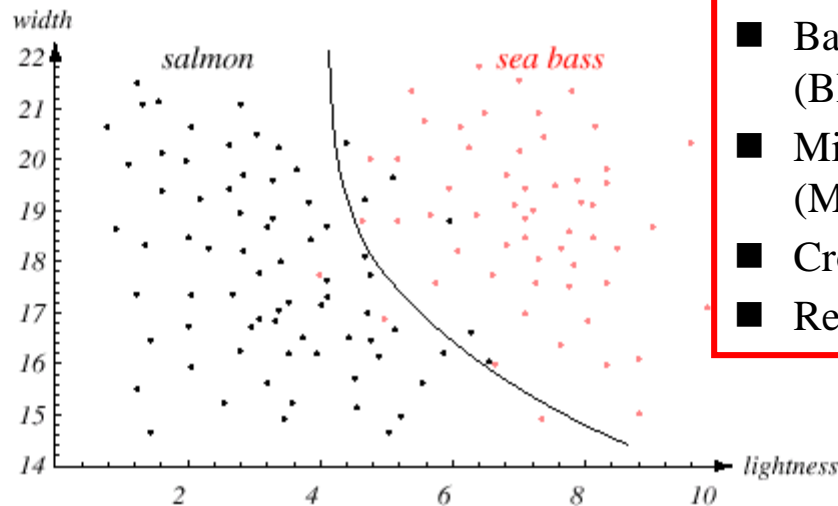


FIGURE 1.5. Overly complex models for the fish will lead to decision boundaries that are complicated. While such a decision may lead to perfect classification of our training samples, it would lead to poor performance on future patterns. The novel test point marked ? is evidently most likely a salmon, whereas the complex decision boundary shown leads it to be classified as a sea bass. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- **Tradeoff between overtraining and generalizability (accuracy in use)**
- **Complex classifiers are not necessarily the best (“Occam’s Razor”)**

How to Select the Right Decision Boundary?

From Duda, Hart and Stork



Model Selection Problem

- Akaike Information Criterion (AIC)
- Bayesian Information Criterion (BIC)
- Minimum Description Length (MDL)
- Cross validation
- Regularization

FIGURE 1.6. The decision boundary shown might represent the optimal tradeoff between performance on the training set and simplicity of classifier, thereby giving the highest accuracy on new patterns. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- **Tradeoff between overtraining and generalizability**
- **How to find the “right” tradeoff?**



IRIS Classification Results using MATLAB APP

Classifier	% Accuracy with 4 Features	% Accuracy with 2 PCA-Reduced Features	% Accuracy with 8 Features*	% Accuracy with 3 PCA-reduced Features*
Decision Tree	94.7%	93.3%	96%	92.7%
Weighted kNN	96.0%	96.7%	95.3%	97.3%
Linear SVM	96.0%	96%	94.7%	98%
Medium Gaussian SVM	96.7%	97.3%	95.3%	96.7%
Bagged Trees	95.3%	92%	96.7%	97.3%

- Add sepal and petal areas and ratios as 4 additional features



Problem 2: Clustering/Categorization - 1

□ **Clustering / Categorization** -- also known as unsupervised classification

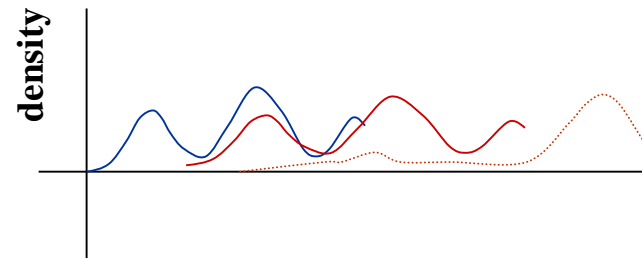
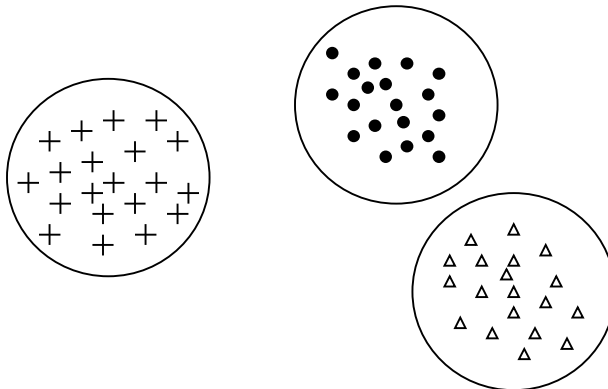
- Explore similarity between the patterns and place similar things in a group



- Derive a statistical model of the process



- Density estimation

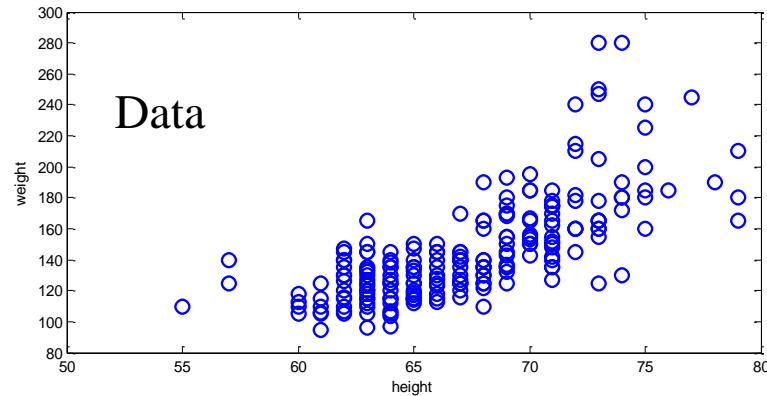


Is there a structure in the data?

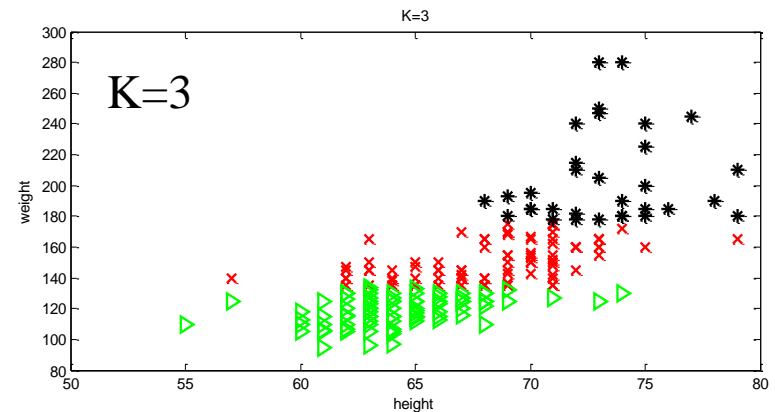
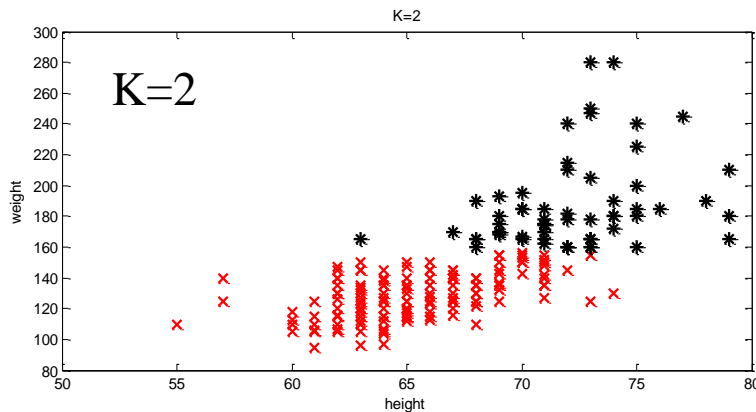


Problem 2: Clustering/Categorization - 2

- K-means clustering of height and weight data



Run `kmeansHeightWeight` in `pmtk3-master\demos` of Murphy to get these

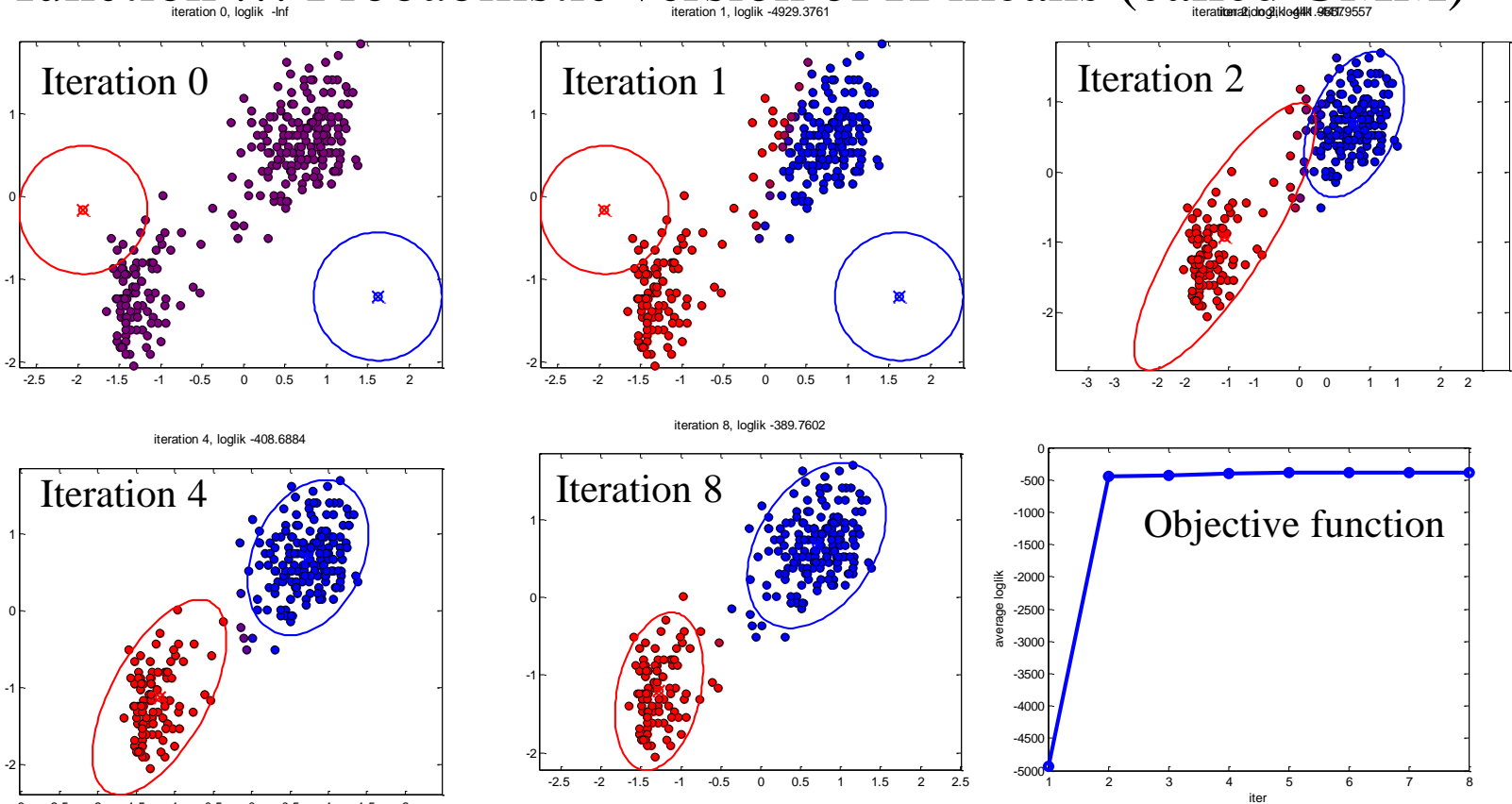


- How to select K? Model selection problem again!



Problem 2: Clustering/Categorization - 3

- How is clustering done? Iteratively by optimizing an objective function ... Probabilistic version of K-means (called GMM)



Old Faithful Data: Run `mixGaussDemoFaithful` in `pmtk3-master\demos` of Murphy to get these

272 observations of waiting time between eruptions and the duration of the eruption for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA.



Problem 3: Function Approximation - 1

□ Function Approximation (Regression)

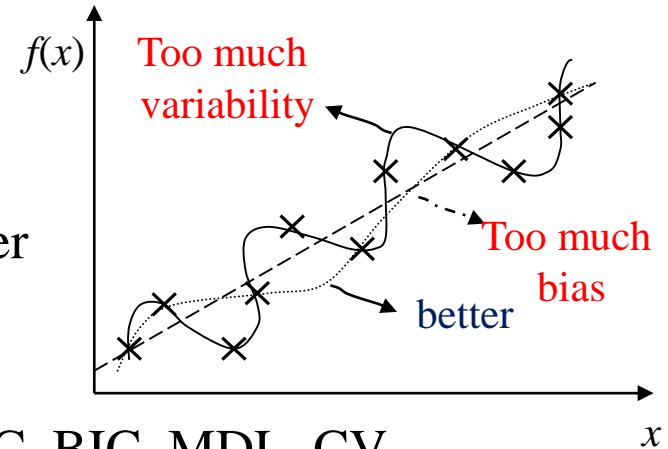
Given $(x_1, z_1), (x_2, z_2) \dots \dots, (x_N, z_N)$, N data points

What is $z \approx f(x)$?

- Very low order models are not good
⇒ too much bias
- Very high order models are not good either
⇒ too much variability

“BIAS-VARIANCE DILEMMA”

- How do we select the “right” model?? AIC, BIC, MDL, CV,
- Applications of function approximation
 - Regression
 - System Identification/Parameter Estimation
 - Neuro-dynamic Programming (NDP)/Approximate DP (ADP)
- What does \approx mean? ⇒ Criterion for optimization of weights



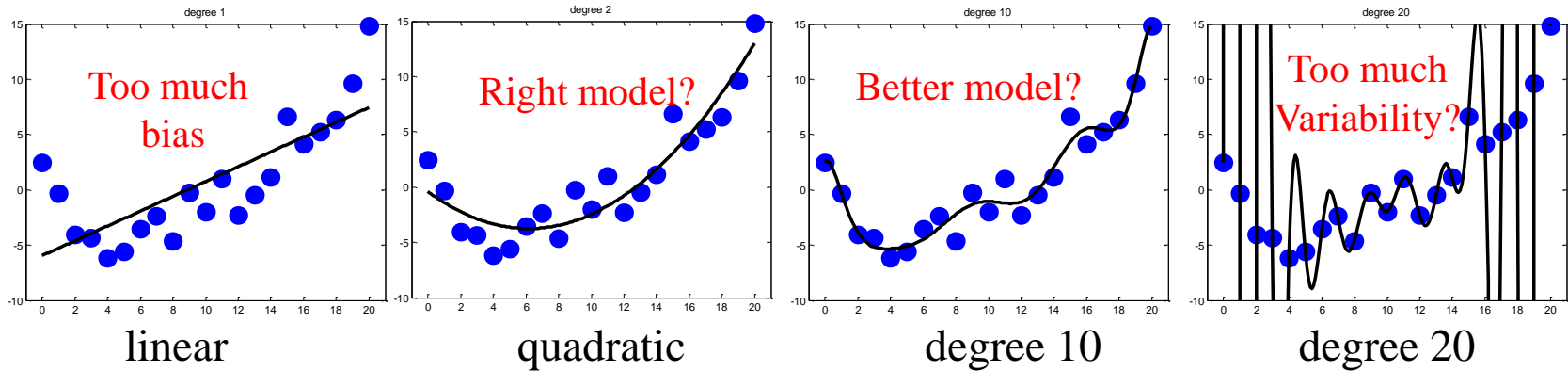
$$\text{MSE: } \frac{1}{N} \sum_{i=1}^N [z_i - f(x_i)]^2$$

This alone is not adequate! Why?

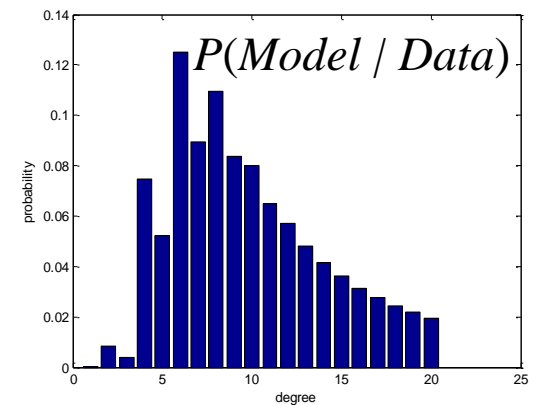
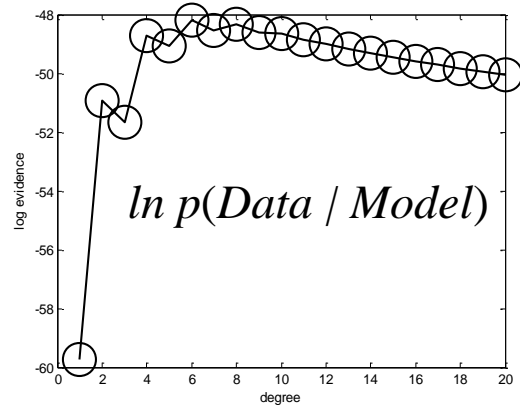
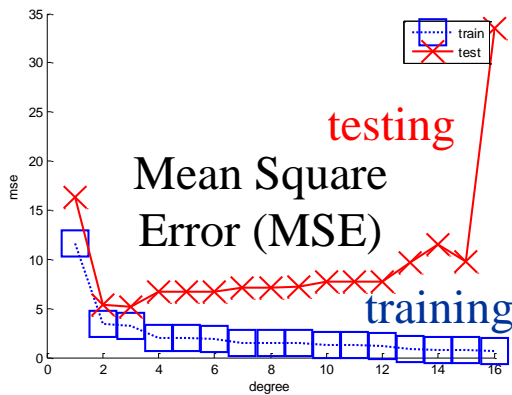


Problem 3: Function Approximation - 2

□ What is the best polynomial for a 21 point noisy data?



Run `linregPolyVsDegree` in `pmtk3-master\demos` of Murphy to get these



Model selection is a difficult problem!



Problem 4: Estimation & Prediction - 1

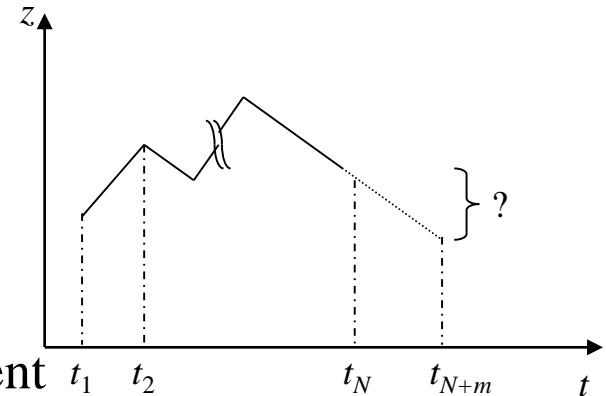
□ Prediction/Forecasting (Regression)

Given $\{z(t_1), z(t_2), \dots, z(t_N)\}$, predict $\hat{z}(t_{N+m})$

□ **Areas:** Estimation theory, Adaptive Filtering

□ Typical Applications

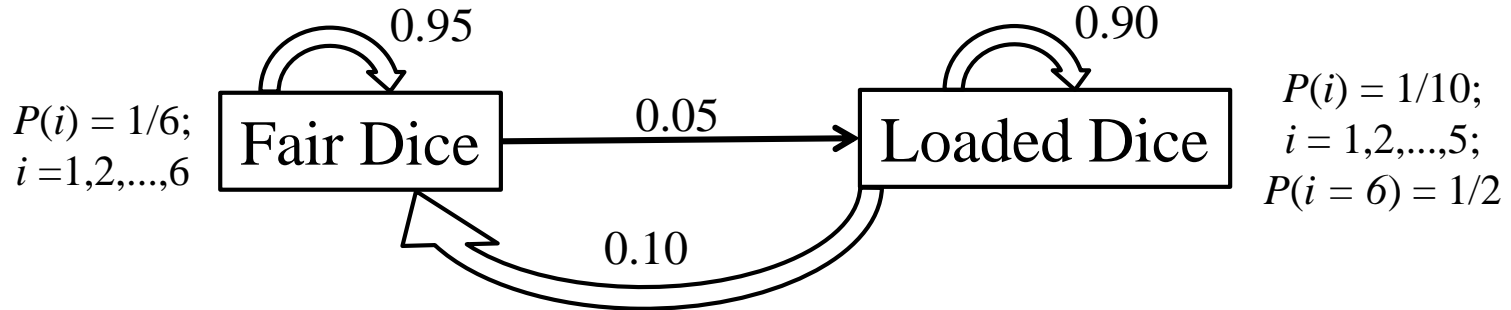
- Stock market prediction
- Weather forecasting
- Prognostics for system health management
- Ad selection systems
- Medicine
- Fraud detection
- Power Demand Prediction



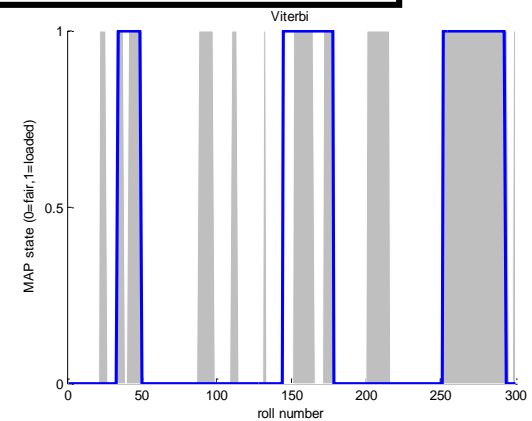
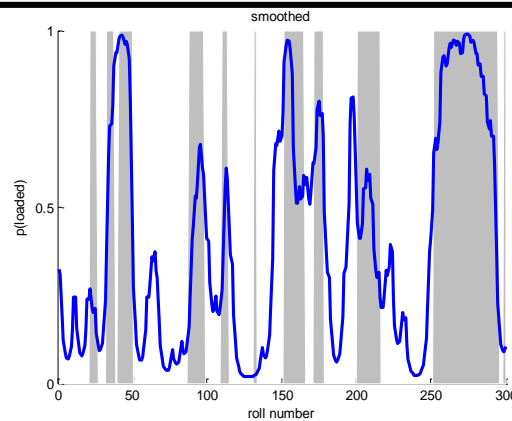
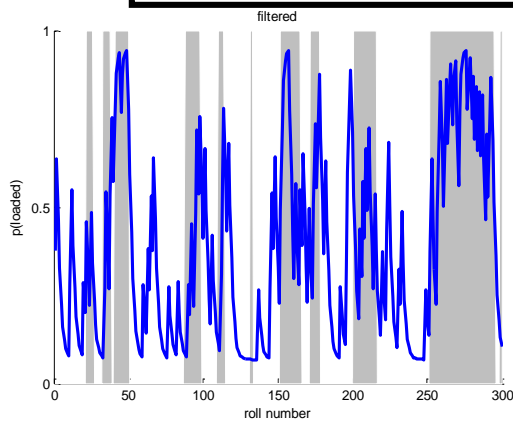


Problem 4: Estimation & Prediction - 2

- How to detect an occasionally dishonest Casino? “Classification”



Run `casinoDemo` in `pmtk3-master\demos` of Murphy to get these



P (Loaded Dice | Data)
(On-line Filtered Estimate)

P (Loaded Dice | Data)
(Off-line Smoothed Estimate)

Decisions: Loaded or not
(Viterbi, MAP Estimate)

- Learn the model based on data & use it to forecast future behavior of Casino



Problem 5: Optimization is Key to Learning

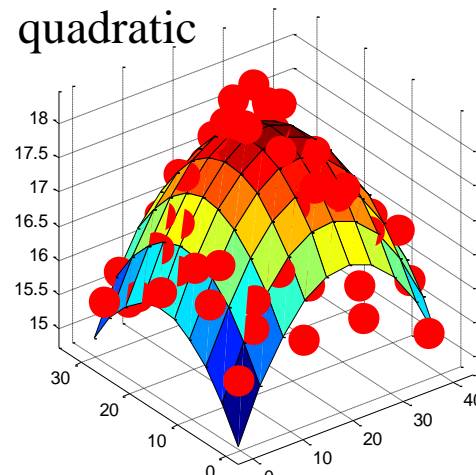
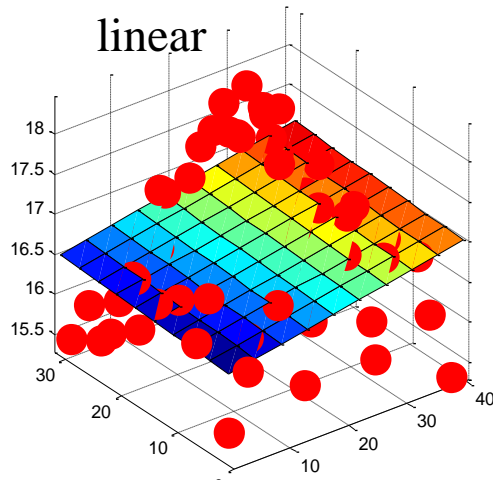
□ Optimization

$$\text{"min } f(\underline{x}) \text{ s.t. } \underline{x} \in \Omega \text{"}$$

$$\Omega \left\{ \begin{array}{l} h(\underline{x}) = \underline{0} \\ g(\underline{x}) \leq \underline{0} \\ \underline{x}_{LB} \leq \underline{x} \leq \underline{x}_{UB} \\ x_i \text{ integer} \end{array} \right.$$

Non-Convex
Non-Differentiable
Combinational (e.g., TSP)

Run surfaceFitDemo in pmtk3-master\demos of Murphy to get these

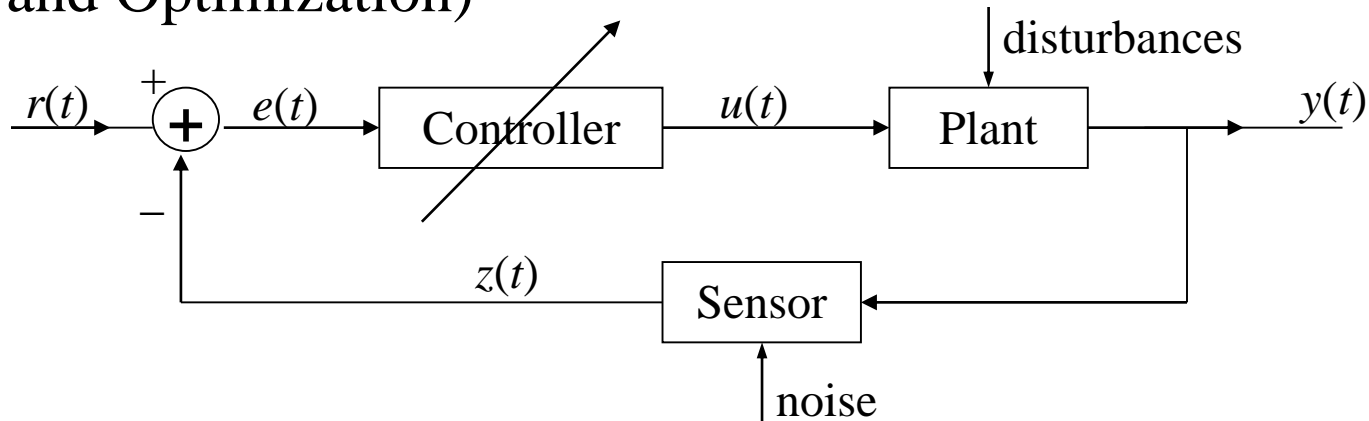


Temperature as a function of (x, y) location in a room

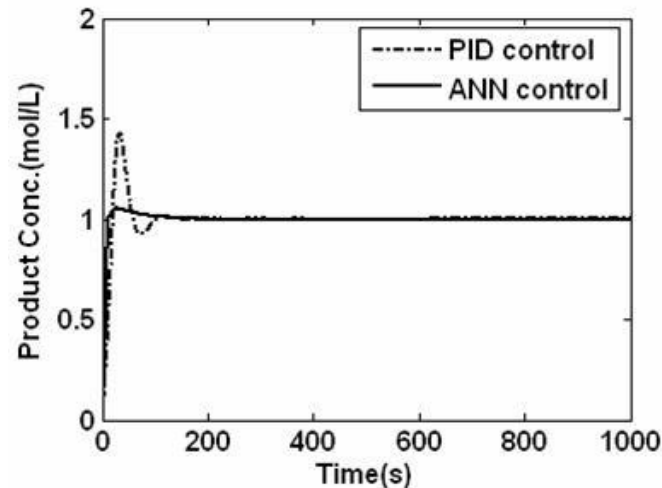


Problem 6: Adaptive Control

- Adaptive (Neuro) Control (Regression, Classification and Optimization)



Design of a controller when have uncertain info. about the system

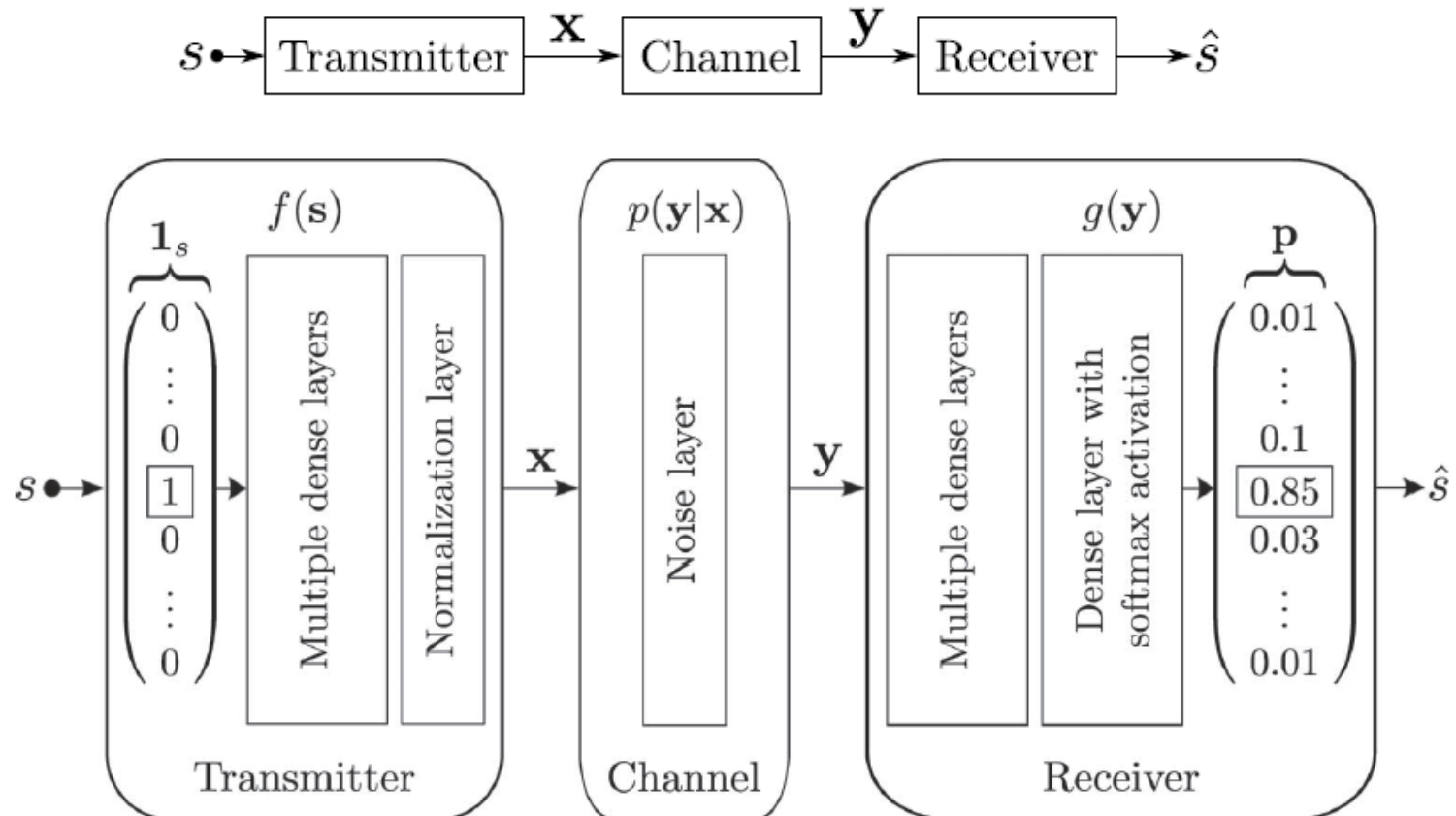


Lots of tuning required to make PID work



Problem 7: Adaptive Communications

- Deep Neural Networks for Physical Layer Design*



* T. O'Shea and J. Hoydis, "An Introduction to Deep Learning for the Physical Layer," *IEEE Trans. on Cognitive Communications and Networking*, Vol. 3, No. 4, pp. 563-575, Dec 2017.



Model Selection Problem is Ubiquitous

- How many clusters in the Data?
- What is the order of the system?
- What is the intrinsic dimension of the system?
- Is this feature relevant?
- How many states in a Hidden Markov Model (HMM)?
- What is the structure of a graphical model?
- How many independent sources in the input?

Bayesian Methods provide a Principled Approach to Model Selection



Bayes' Theorem and Inference

□ Basic Axioms of Probability

- Probability of event A , $P(A) \in [0, 1]$
- $P(A) = 1 \Leftrightarrow A$ is certain
- $P(A \cup B) = P(A \text{ or } B) = P(A) + P(B) - P(A \cap B)$
- Bayes' theorem

- $P(AB) = P(A | B)P(B) = P(B | A)P(A)$

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

- Interested in A
- Begin with *a priori* probability $P(A)$ for our belief about A
- Observe B
- Bayes' theorem provides the revised belief about A , that is, the posterior probability $P(A/B)$
- Likelihood of A : The quantity $P(B/A)$, as a function of varying A for a fixed B
- *posterior* \propto *prior* \times *likelihood*
 $\Rightarrow P(A/B) \propto P(A) \cdot P(B/A)$

We are inferring A given data B

• Sum Rule (Total Probability Theorem)

$$P(A) = \sum_B P(A, B)$$

• Product Rule

$$P(A, B) = P(B | A)P(A) = P(A | B)P(B)$$

• $P \sim$ Probability

• $p \sim$ Probability Density Function (pdf)



Bayes' Theorem and Learning

- Data D ; Model Parameters, \underline{w} ; model set $M \in \{1, 2, \dots, m, \dots, M\}$

$$p(\underline{w} | D, m) = \frac{p(D | \underline{w}, m) p(\underline{w} | m)}{p(D | m)}$$

$p(D | \underline{w}, m)$ = Likelihood of \underline{w} under model m

$p(\underline{w} | m)$ = prior of \underline{w} under model m

$p(\underline{w} | D, m)$ = posterior of \underline{w} given data D under model m

- Prediction

$$p(\underline{x} | D, m) = \int p(\underline{x}, \underline{w} | D, m) d\underline{w} = \int p(\underline{x} | \underline{w}, m) p(\underline{w} | D, m) d\underline{w}$$

- Model Comparison

$$P(m | D) = \frac{p(D | m)P(m)}{\sum_{k=1}^M p(D | k)P(k)}; p(D | m) = \int p(D | \underline{w}, m) p(\underline{w} | m) d\underline{w} \sim \text{Model Evidence}$$

Interpretations of Model Evidence:

- Pdf that randomly selected model parameters from the prior generate D
- $-\log_2 p(D/m)$ is the *number of bits of surprise* at observing D under model m
- Bayesian model averaging Computationally intractable in general

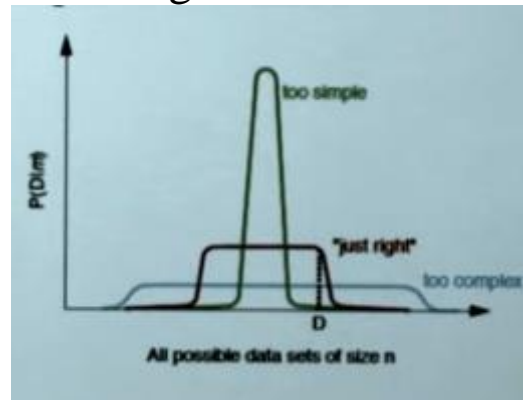


Bayesian Recursive Learning and Inference

- ❑ Learning and Inference are not different in a Bayesian approach
- ❑ Recursion

$$p(\underline{w} | D, \underline{x}, m) = \frac{p(\underline{w}, D, \underline{x} | m)}{p(D, \underline{x} | m)} = \frac{p(\underline{x} | \underline{w}, m) p(\underline{w} | D, m)}{\int p(\underline{x} | \underline{w}, m) p(\underline{w} | D, m) d\underline{w}}$$

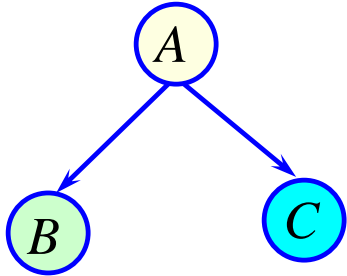
- ❑ Provides a principled method^w for Model Selection
 - Very Simple model structures are unlikely to generate observed data (“large bias”)
 - Very Complex model structures are unlikely to generate the particular data set observed because they can generate many possible data sets (“large variance”)
 - The Right model is one with large model evidence for the observed data set D Large $p(D|m)$





Graphical Modeling of Statistical Dependencies – 1

Example 1 (Common Cause, Fork)



- Naïve Bayes
- iid observations
- prediction (sufficient statistics)

When A is not observed, B & C are **dependent**.
 When A is observed, B & C are **conditionally independent!**

“tail-to-tail” dependency

$$B \perp C | A$$

"B is independent of C given A"

Why?

$$P(B, C) = \sum_A P(B | A)P(C | A)P(A) \neq P(B)P(C)$$

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

Discuss chain rule:

$$P(ABC) = P(A)P(B | A)P(C | A, B)$$

$$P(ABC) = P(A)P(B | A)P(C | A)$$

$$= \frac{P(AB)P(AC)}{P(A)}$$

Exploiting graph structure

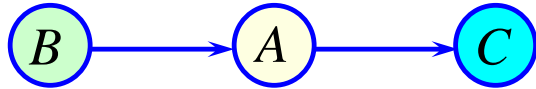
Examples: 1. Smoking (A) causes Bronchitis (B) and Cancer (C).

2. Learning the probability of heads $P(A=1)$, given the outcomes of a number of coin tosses (in this case, two (B & C))



Graphical Modeling of Statistical Dependencies – 2

Example 2 (Indirect Cause, Chain)



“head-to-tail” dependency

$$P(ABC) = P(B)P(A|B)P(C|A)$$

$$\text{Also, } P(ABC) = \frac{P(AB)P(CA)}{P(A)}$$

{AB}, {CA} “cliques”

A “separator”

$P(AB)$, $P(CA)$ are “clique potentials”

$P(A)$ “separator potential”

When A is not observed, B & C are **dependent**.

When A is observed, B & C are **conditionally independent!**

$$B \perp C \mid A$$

"B is independent of C given A"

Why?

$$P(B, C) = P(B) \sum_A P(A|B)P(C|A) \neq P(B)P(C)$$

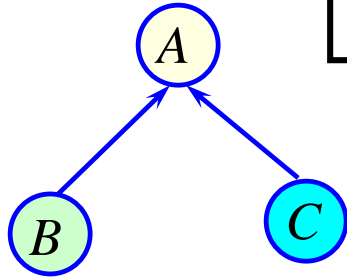
$$P(B, C \mid A) = \frac{P(A, B, C)}{P(A)} = \frac{P(A|B)P(B)P(C|A)}{P(A)} = P(B|A)P(C|A)$$

Example: B: Clouds; A: Rain; C: Grass is wet



Graphical Modeling of Statistical Dependencies – 3

- Example 3 (Common Effect, Collider)



“head-to-head” dependency

$$P(ABC) = P(B)P(C)P(A|B,C)$$

When A is not observed, B & C are **independent!!**
When A is observed, B & C are **conditionally dependent!!**

$$B \perp C | \emptyset$$

"B is independent of C given no evidence"

Why?

$$P(B, C) = \sum_A P(B)P(C)P(A|B, C) = P(B)P(C)$$

They are **not** independent given A

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

When A is not observed, A blocks the path from B to C. However, when A is observed, it **unblocks** the path from B to C \Rightarrow *they become dependent*

Example: A: Engine Does Not Start; B: Battery; C: Fuel

Noisy (“Soft”) XOR

B	C	$P(A=1 B,C)$	$P(A=0 B,C)$
0	0	0.10	0.90
0	1	0.99	0.01
1	0	0.80	0.20
1	1	0.25	0.75

$$P(B=1)=0.65$$

$$P(C=1)=0.77$$

Suppose we observed $A = 0$. What is $P(B = 1 | A = 0)$?

$$P(B = 1 | A = 0) = \frac{P(A = 0 | B = 1)P(B = 1)}{P(A = 0)} = \frac{P(A = 0 | B = 1)P(B = 1)}{P(A = 0 | B = 1)P(B = 1) + P(A = 0 | B = 0)P(B = 0)}$$

$$P(A = 0 | B = 1) = \sum_{C=0}^1 P(A = 0, C | B = 1) = \sum_{C=0}^1 P(A = 0 | B = 1, C)P(C | B = 1)$$

$$= \sum_{C=0}^1 P(A = 0 | B = 1, C)P(C) = 0.20 * 0.23 + 0.75 * 0.77 = 0.6235$$

$$\text{Similarly, } P(A = 0 | B = 0) = \sum_{C=0}^1 P(A = 0 | B = 0, C)P(C) = 0.90 * 0.23 + 0.01 * 0.77 = 0.2147$$

$$P(B = 1 | A = 0) = \frac{0.6235 * 0.65}{0.6235 * 0.65 + 0.2147 * 0.35} = 0.8436$$

$$P(B = 0 | A = 0) = 1 - P(B = 1 | A = 0) = 0.1564$$



Exploiting Dependency Structure: Graphical Models- 1

□ Chain rule (general dependency)

$$p(x_1, x_2, x_3, \dots, x_{p-1}, x_p) = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) \dots p(x_p | x_1, x_2, \dots, x_{p-1})$$

- Storage complexity $O(K^p)$, assuming each variable has K states

□ Conditional Independence

- Naïve Bayes

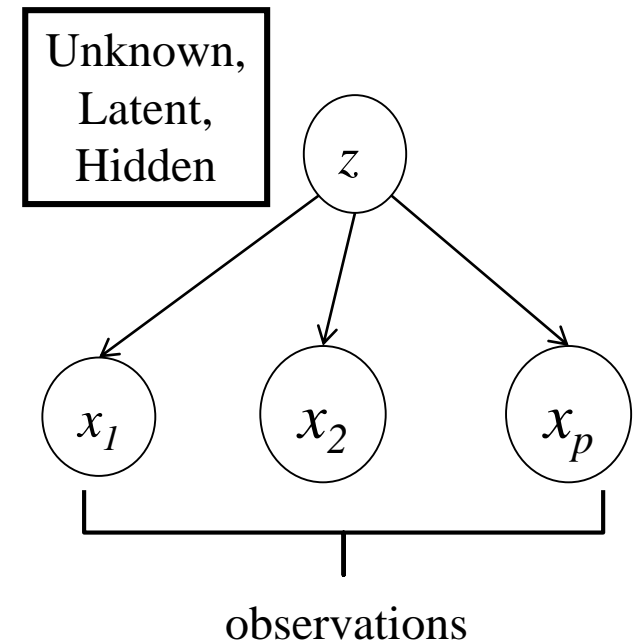
$$p(\underline{x}, z) = p(z) \left(\prod_{i=1}^p p(x_i | z) \right)$$

- Storage complexity $[(K-1)(pC+1)]$

$$x_i \in \{1, 2, \dots, K\}$$

$$z \in \{1, 2, \dots, C\}$$

- Classifiers based on Naïve Bayes can work quite well

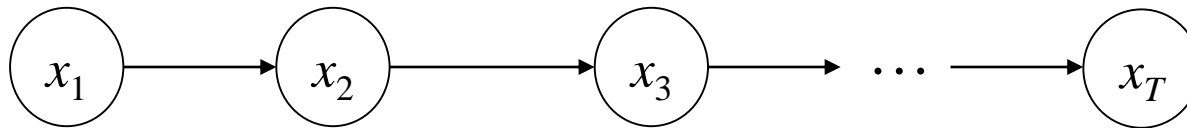




Exploiting Dependency Structure: Graphical Models - 2

□ Conditional Independence (continued)

- Markov chain (used to model sequential first order dependencies)



bigram

$$p(x_{1:T}) = p(x_1) \prod_{t=2}^T p(x_t | x_{t-1})$$

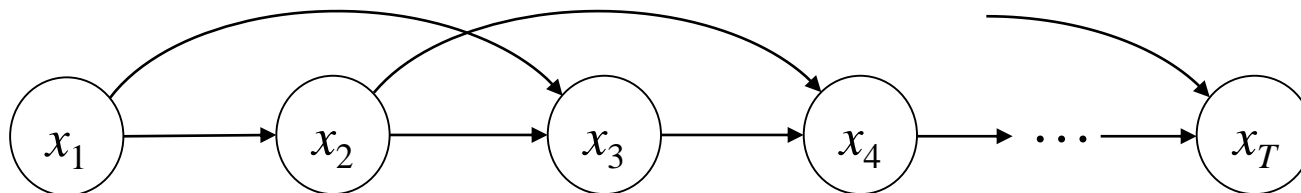
- Homogeneous: K^2
- Non-homogeneous: TK^2

Applications:

- Sentence completion
- Data compression
- Text classification
- Automatic essay writing
- Google's page rank

- n^{th} order Markov chain ((**n+1**)-gram models)

$$p(x_{1:T}) = p(x_1, x_2, \dots, x_n) \prod_{t=n+1}^T p(x_t | x_{t-1}, x_{t-2}, \dots, x_{t-n}); T > n$$



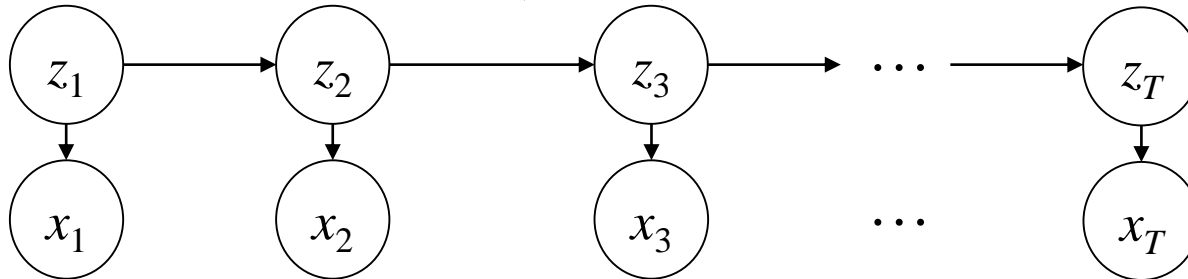
$n = 2$
trigram



Exploiting Dependency Structure: Graphical Models - 3

□ Conditional Independence (continued)

- Hidden Markov model (HMM)



$$p(z_{1:T}, x_{1:T}) = p(z_{1:T}) p(x_{1:T} | z_{1:T}) = p(z_1) \cdot \left(\prod_{t=2}^T p(z_t | z_{t-1}) \right) \cdot \left(\prod_{t=1}^T p(x_t | z_t) \right)$$

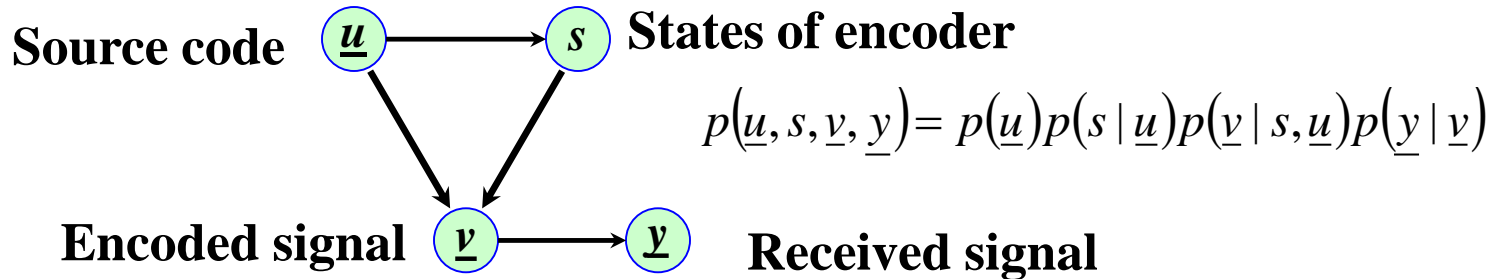
- **Applications:** speech recognition, activity recognition, gene finding, protein sequence alignment, diagnosis,...
- **Algorithms:** Forwards-backwards algorithms, Viterbi decoding, Baum-Welch algorithm for learning HMM parameters from data
- **Generalizations:** Higher order, Factorial, Coupled, semi-Markov, Hierarchical, I/O, Auto-regressive, Buried, state space, Dynamic Bayesian networks (DBNs),...



Three Broad Classes of Graphical Models - 1

Bayesian Networks

- Represented in terms of directed acyclic graphs (DAGs)



- In general,

$$\mathbf{z} = [z_1 \ z_2 \ \dots \ z_N]$$

$$P(z_k | a_k) \quad a_k = \text{parents of } z_k = pa(z_k)$$

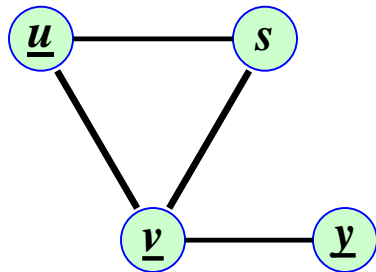
$$P(\mathbf{z}) = \prod_{k=1}^N P(z_k | pa(z_k)) = \prod_{k=1}^N P(z_k | a_k)$$

- Exploit graph properties
 - d-separation, moralization, cliques,...
- Algorithms
 - Variable elimination
 - Junction tree
 - Variational inference
 - Belief propagation (BP)
 - Expectation propagation (EP)
 - Markov Chain Monte Carlo



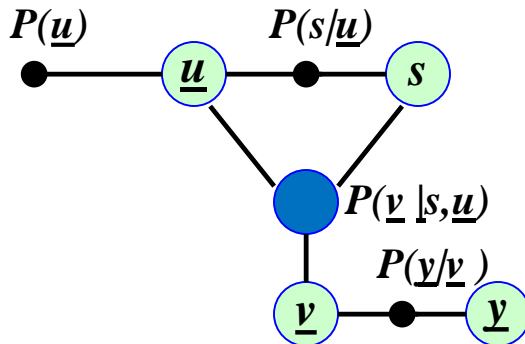
Three Broad Classes of Graphical Models - 2

- Markov random fields (MRF) and conditional MRFs...
Undirected graphical model (UGM) or Markov network



□ Factor Graphs

- Useful when dealing with large graphs
- Both BNs and UGMs can be converted into factor graphs



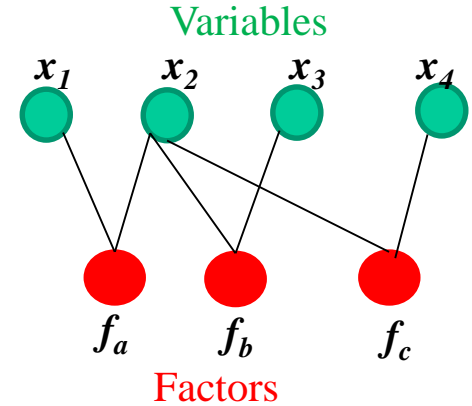
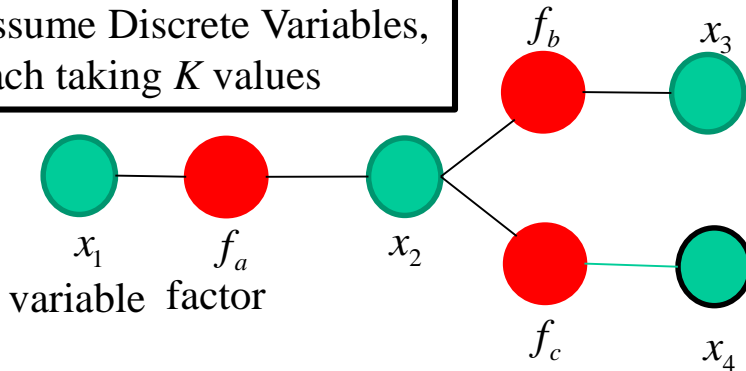
$$p(\underline{u}, s, \underline{v}, \underline{y}) = \underbrace{p(\underline{u})}_{\text{factor}} \underbrace{p(s|\underline{u})}_{\text{factor}} \underbrace{p(\underline{v}|s,\underline{u})}_{\text{factor}} \underbrace{p(\underline{y}|\underline{v})}_{\text{factor}}$$

Propagate beliefs between
variables \rightarrow factors &
factors \rightarrow variables



Efficient Inference via Message Passing - 1

Assume Discrete Variables,
Each taking K values



Corresponding Bipartite Graph

$$P(x_1, x_2, x_3, x_4) = f_a(x_1, x_2) f_b(x_2, x_3) f_c(x_2, x_4)$$

Suppose want marginal distribution $P(x_2)$. By sum formula or Total Probability Theorem

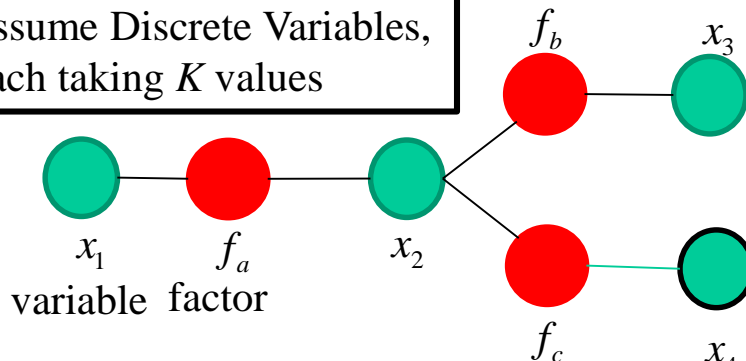
$$\begin{aligned}
 P(x_2) &= \sum_{x_1} \sum_{x_3} \sum_{x_4} f_a(x_1, x_2) f_b(x_2, x_3) f_c(x_2, x_4) \\
 &= \underbrace{\left(\sum_{x_1} f_a(x_1, x_2) \right)}_{\mu_{f_a \rightarrow x_2}(x_2)} \underbrace{\left(\sum_{x_3} f_b(x_2, x_3) \right)}_{\mu_{f_b \rightarrow x_2}(x_2)} \underbrace{\left(\sum_{x_4} f_c(x_2, x_4) \right)}_{\mu_{f_c \rightarrow x_2}(x_2)} = \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_b \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2)
 \end{aligned}$$

Note : $2K^3$ multiplications and K^3 additions versus 2 multiplications and $3K$ additions for each x_2 ! In fact, can compute any marginal distribution by exchanging messages from variables to factors and factors to variables.



Efficient Inference via Message Passing - 2

Assume Discrete Variables,
Each taking K values



To compute all marginals, form an arbitrary root node (say x_1), propagate and store messages from the leaves to the root, and then propagate and store messages from the root to the leaves. Marginals are products of the factor messages received at each variable nodes. Normalize marginals Exact for trees. For arbitrary graphs, iterate (approximate belief propagation (ABP)). In fact, you can pass messages asynchronously from nodes \rightarrow factors and factors \rightarrow nodes until convergence.

$$P(x_1) = \mu_{f_a \rightarrow x_1}(x_1); P(x_3) = \mu_{f_b \rightarrow x_3}(x_3); P(x_4) = \mu_{f_c \rightarrow x_4}(x_4)$$

$$P(x_2) = \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_b \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2)$$

Leaves \rightarrow Root

$$\mu_{x_3 \rightarrow f_b}(x_3) = 1 \forall x_3$$

$$\mu_{f_b \rightarrow x_2}(x_2) = \sum_{x_3} f_b(x_2, x_3) \mu_{x_3 \rightarrow f_b}(x_3)$$

$$\mu_{x_4 \rightarrow f_c}(x_4) = 1 \forall x_4$$

$$\mu_{f_c \rightarrow x_2}(x_2) = \sum_{x_4} f_c(x_2, x_4) \mu_{x_4 \rightarrow f_c}(x_4)$$

$$\mu_{x_2 \rightarrow f_a}(x_2) = \mu_{f_b \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2)$$

$$\mu_{f_a \rightarrow x_1}(x_1) = \sum_{x_2} f_a(x_1, x_2) \mu_{x_2 \rightarrow f_a}(x_2)$$

Root \rightarrow Leaves

$$\mu_{x_1 \rightarrow f_a}(x_1) = 1 \forall x_1$$

$$\mu_{f_a \rightarrow x_2}(x_2) = \sum_{x_1} f_a(x_1, x_2) \mu_{x_1 \rightarrow f_a}(x_1)$$

$$\mu_{x_2 \rightarrow f_b}(x_2) = \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2)$$

$$\mu_{x_2 \rightarrow f_c}(x_2) = \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_b \rightarrow x_2}(x_2)$$

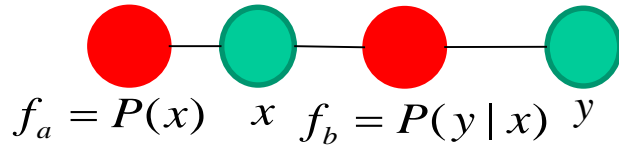
$$\mu_{f_b \rightarrow x_3}(x_3) = \sum_{x_2} f_b(x_2, x_3) \mu_{x_2 \rightarrow f_b}(x_2)$$

$$\mu_{f_c \rightarrow x_4}(x_4) = \sum_{x_2} f_c(x_2, x_4) \mu_{x_2 \rightarrow f_c}(x_2)$$

Evidence is easy to include



Bayes Rule as Message Passing



$$P(x) = \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}; x \in \{1, 2\}$$

$$P(y | x) = \begin{bmatrix} 0.05 & 0.65 & 0.30 \\ 0.8 & 0.10 & 0.10 \end{bmatrix}; y \in \{1, 2, 3\}$$

Want $P(y)$ and $P(x | y = 1)$

Marginal $P(y)$

$$\mu_{f_a \rightarrow x}(x) = P(x) \Rightarrow \mu_{f_a \rightarrow x}(x = 1) = 0.8; \mu_{f_a \rightarrow x}(x = 2) = 0.2 \quad \text{Bishop's "Who done it" Example}$$

$$\mu_{x \rightarrow f_b}(x) = \mu_{f_a \rightarrow x}(x)$$

$$\mu_{f_b \rightarrow y}(y) = \sum_x f_b(x, y) \mu_{f_a \rightarrow x}(x) \Rightarrow \mu_{f_b \rightarrow y}(y = 1) = 0.2; \mu_{f_b \rightarrow y}(y = 2) = 0.54; \mu_{f_b \rightarrow y}(y = 3) = 0.26$$

$$P(y) = \mu_{f_b \rightarrow y}(y)$$

Evidence, $y = 1$. What is $P(x | y = 1)$?

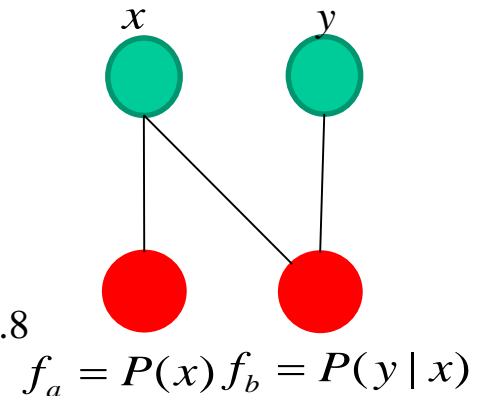
Treat x as the root.

$$\mu_{y \rightarrow f_b}(y = 1) = 1; \mu_{y \rightarrow f_b}(y = 2) = 0; \mu_{y \rightarrow f_b}(y = 3) = 0$$

$$\mu_{f_b \rightarrow x}(x) = \sum_{y=1}^3 f_b(x, y) \mu_{y \rightarrow f_b}(y) \Rightarrow \mu_{f_b \rightarrow x}(x = 1) = 0.05; \mu_{f_b \rightarrow x}(x = 2) = 0.8$$

$$\mu_{f_a \rightarrow x}(x) = P(x) \Rightarrow \mu_{f_a \rightarrow x}(x = 1) = 0.8; \mu_{f_a \rightarrow x}(x = 2) = 0.2$$

$$P(x | y = 1) = \mu_{f_a \rightarrow x}(x) \mu_{f_b \rightarrow x}(x) \Rightarrow \begin{bmatrix} \mu_{f_a \rightarrow x}(1) \mu_{f_b \rightarrow x}(1) \\ \mu_{f_a \rightarrow x}(2) \mu_{f_b \rightarrow x}(2) \end{bmatrix} = \begin{bmatrix} 0.04 \\ 0.16 \end{bmatrix} \Rightarrow \text{Normalize: } \begin{bmatrix} 0.20 \\ 0.80 \end{bmatrix}$$





Sum-Product Belief Propagation Algorithm

- **Factors → Variables Messages (SUM)**

messages sent by a factor node to a variable node involves multiplying all the incoming messages (except variable node x) with the factor and summing over all the variables except x

$$\mu_{f_a \rightarrow x_1}(x_1) = \sum_{x_2} f_a(x_1, x_2) \mu_{x_2 \rightarrow f_a}(x_2); \mu_{f_a \rightarrow x_2}(x_2) = \sum_{x_1} f_a(x_1, x_2) \mu_{x_1 \rightarrow f_a}(x_1)$$

$$\mu_{f_b \rightarrow x_2}(x_2) = \sum_{x_3} f_b(x_2, x_3) \mu_{x_3 \rightarrow f_b}(x_3); \mu_{f_b \rightarrow x_3}(x_3) = \sum_{x_2} f_b(x_2, x_3) \mu_{x_2 \rightarrow f_b}(x_2)$$

$$\mu_{f_c \rightarrow x_2}(x_2) = \sum_{x_4} f_c(x_2, x_4) \mu_{x_4 \rightarrow f_c}(x_4); \mu_{f_c \rightarrow x_4}(x_4) = \sum_{x_2} f_c(x_2, x_4) \mu_{x_2 \rightarrow f_c}(x_2)$$

- **Variables → Factors Messages (PRODUCT)**

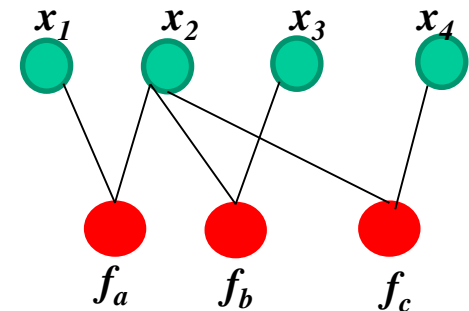
Message sent by a variable node to a factor node is the product of all the incoming messages along all of the other links (factors)

$$\mu_{x_1 \rightarrow f_a}(x_1) = 1; \mu_{x_3 \rightarrow f_b}(x_3) = 1; \mu_{x_4 \rightarrow f_c}(x_4) = 1$$

$$\mu_{x_2 \rightarrow f_a}(x_2) = \mu_{f_b \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2);$$

$$\mu_{x_2 \rightarrow f_b}(x_2) = \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2);$$

$$\mu_{x_2 \rightarrow f_c}(x_2) = \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_b \rightarrow x_2}(x_2)$$





Exponential Family

□ Exponential family:

Gaussian, Poisson, Gamma, Binomial, Exponential, Beta, Weibull,...

$$p(\underline{x} | \underline{w})^\dagger = \frac{1}{Z(\underline{w})} h(\underline{x}) \exp\{\underline{w}^T \underline{T}(\underline{x})\} = h(\underline{x}) \exp\{\underline{w}^T \underline{T}(\underline{x}) - A(\underline{w})\}$$

$$Z(\underline{w}) = \int_{\underline{x}} h(\underline{x}) \exp\{\underline{w}^T \underline{T}(\underline{x})\} d\underline{x} \quad \dots \text{ called partition function}$$

$A(\underline{w}) = \ln Z(\underline{w})$... log partition function or cumulant function

Gradients of $A(\underline{w})$ give moments of density

$\underline{w} \sim$ canonical parameters

$\underline{T}(\underline{x}) =$ "sufficient statistics"

$h(\underline{x}) =$ scaling constant

† most often $\underline{\theta}$ is used to denote model parameters. Since Neural Networks use \underline{w} for weights, we use \underline{w} here.



Exponential Family and Conjugate Priors

□ Conjugate Prior

$p(\underline{w}) = g(\underline{\eta}, \underline{v}) \exp\{\underline{w}^T \underline{v} - \eta A(\underline{w})\}$; $\underline{\eta}, \underline{v}$ are hyper parameters

Data $D = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$

$p(\underline{x}_n | \underline{w}) = h(\underline{x}_n) \exp\{\underline{w}^T \underline{T}(\underline{x}_n) - A(\underline{w})\}$

$p(\underline{w} | D) = g(\underline{\eta} + N, \underline{v} + \sum_{n=1}^N \underline{T}(\underline{x}_n)) \exp\{\underline{w}^T [\underline{v} + \sum_{n=1}^N \underline{T}(\underline{x}_n)] - (\eta + N)A(\underline{w})\}$

Posterior belongs to the exponential family

□ Examples of Likelihoods and Conjugate Priors

- (Bernoulli or Binomial or Geometric) –Beta
- Scalar Normal-(Normal, Inverse-Gamma)
- Multivariate Normal-(Normal, Inverse Wishart)
- (Poisson, Pareto, Exponential)-Gamma



Bernoulli Distribution

□ Bernoulli

$$p(x) = p^x (1-p)^{1-x} = \exp[x \ln p + (1-x) \ln(1-p)] = \exp\left[x \ln \frac{p}{1-p} + \ln(1-p)\right]$$

$$\text{Let } T(x) = x; w = \ln \frac{p}{1-p} \Rightarrow p = \frac{1}{1+e^{-w}} = \sigma(w)$$

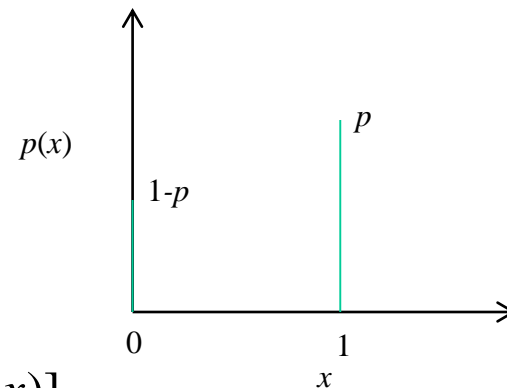
$$\Rightarrow 1-p = \frac{1}{1+e^w} = \sigma(-w) = 1 - \sigma(w)$$

$$\text{so, } p(x) = \exp\left[x \ln \frac{p}{1-p} - \ln(1+e^w)\right]$$

$$A(w) = \ln(1+e^w)$$

$$\frac{dA(w)}{dw} = \frac{e^w}{1+e^w} = \frac{1}{1+e^{-w}} = p = E(x) = E[T(x)]$$

$$\frac{d^2 A(w)}{dw^2} = p(1-p) = \text{Var}(x) = \text{Var}[T(x)]$$



Binomial Likelihood -Beta Prior for Learning p

- Multiple trials, each with $x_i \in \{0,1\}$

$$s_n = \sum_{i=1}^n x_i; s_{n+1} = s_n + x_{n+1}; s_0 = 0$$

$$\text{Beta Function: } B(\alpha, \beta) = \int_0^1 q^{\alpha-1} (1-q)^{\beta-1} dq = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} = \frac{(\alpha-1)!(\beta-1)!}{(\alpha+\beta-1)!}$$

$$p(s_n | p) = \binom{n}{s_n} p^{s_n} (1-p)^{n-s_n} = \frac{p^{s_n} (1-p)^{n-s_n}}{(n+1)B(s_n+1, n+1-s_n)}$$

$$p(p) = \frac{p^{\alpha-1} (1-p)^{\beta-1}}{B(\alpha, \beta)} = \text{Beta}(\alpha, \beta)$$

$$p(p | s_n) = \frac{p^{s_n+\alpha-1} (1-p)^{n-s_n+\beta-1}}{\int_0^1 q^{s_n+\alpha-1} (1-q)^{n-s_n+\beta-1} dq} = \frac{p^{s_n+\alpha-1} (1-p)^{n-s_n+\beta-1}}{B(s_n+\alpha, n-s_n+\beta)} = \text{Beta}(s_n+\alpha, n-s_n+\beta)$$

$$\text{Estimate: } \hat{p}_n = \frac{s_n + \alpha}{n + \alpha + \beta}$$

$$\text{Recursion: } \hat{p}_{n+1} = \frac{s_{n+1} + \alpha}{n + \alpha + \beta + 1} = \frac{s_n + x_{n+1} + \alpha}{n + \alpha + \beta + 1} = \frac{n + \alpha + \beta}{n + \alpha + \beta + 1} \hat{p}_n + \frac{1}{n + \alpha + \beta + 1} x_{n+1}$$

$$\text{So, } \hat{p}_{n+1} = \hat{p}_n + \frac{1}{n + \alpha + \beta + 1} (x_{n+1} - \hat{p}_n)$$

Similar to Least Squares and Kalman Filter



Univariate Gaussian (Normal) Distribution

□ Univariate Gaussian Distribution

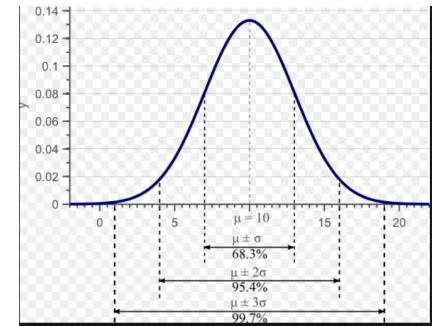
$$p(\underline{x} | \underline{w}) = \frac{1}{Z(\underline{w})} h(\underline{x}) \exp\{\underline{w}^T \underline{T}(\underline{x})\} = h(\underline{x}) \exp\{\underline{w}^T \underline{T}(\underline{x}) - A(\underline{w})\}$$

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

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

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

$$= \underbrace{1}_{h(x)} \exp\left\{\underbrace{\begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix}}_{\underline{w}}^T \underbrace{\begin{pmatrix} x \\ x^2 \end{pmatrix}}_{\underline{T}(x)} - \underbrace{\left(\frac{\mu^2}{2\sigma^2} + \ln \sigma + \frac{1}{2} \ln(2\pi)\right)}_{A(\underline{w})}\right\}$$



μ : mean
 σ^2 : variance
 $1/\sigma^2$: “information”



Log Partition Function: UVN

$$\underline{T}(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}; \underline{w} = \begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\begin{aligned} A(\underline{w}) &= \frac{\mu^2}{2\sigma^2} + \ln \sigma + \frac{1}{2} \ln(2\pi) \\ &= -\frac{w_1^2}{4w_2} - \frac{1}{2} \ln(-2w_2) + \frac{1}{2} \ln(2\pi) \end{aligned}$$

$$\begin{aligned} \frac{\partial A(\underline{w})}{\partial w_1} &= -\frac{w_1}{2w_2} = -\frac{\mu}{\sigma^2} \cdot (-\sigma^2) = \mu = E(x) \\ \frac{\partial A(\underline{w})}{\partial w_2} &= \frac{w_1^2}{4w_2^2} - \frac{1}{2} \cdot \frac{1}{(-2w_2)} \cdot (-2) \\ &= \frac{w_1^2}{4w_2^2} - \frac{1}{2w_2} \\ &= \mu^2 + \sigma^2 = E(x^2) \end{aligned}$$

$$\nabla_{\underline{w}} A(\underline{w}) = E[\underline{T}(x)]$$

$$\frac{\partial^2 A(\underline{w})}{\partial w_1^2} = -\frac{1}{2w_2} = \sigma^2 = \text{var}(x)$$

$$\frac{\partial^2 A(\underline{w})}{\partial w_1 \partial w_2} = \frac{\partial^2 A(\underline{w})}{\partial w_2 \partial w_1} = \frac{w_1}{2w_2^2} = \frac{1}{2} \cdot \frac{\mu}{\sigma^2} \cdot 4\sigma^4 = 2\mu\sigma^2 = E\{(x-\mu)(x^2 - \mu^2 - \sigma^2)\}$$

$$\frac{\partial^2 A(\underline{w})}{\partial w_2^2} = -\frac{w_1^2}{2w_2^3} + \frac{1}{2w_2^2} = \frac{1}{2} \cdot \frac{\mu^2}{\sigma^4} \cdot 8\sigma^6 + 2\sigma^4 = 4\mu^2\sigma^2 + 2\sigma^4 = \text{var}(x^2)$$

$$\nabla_{\underline{w}}^2 A(\underline{w}) = \begin{bmatrix} \sigma^2 & 2\mu\sigma^2 \\ 2\mu\sigma^2 & 4\mu^2\sigma^2 + 2\sigma^4 \end{bmatrix} = \text{Cov}(\underline{T}(x))$$

$$\begin{aligned} E(x^3) &= \mu^3 + 3\mu\sigma^2 \\ E(x^4) &= \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4 \end{aligned}$$



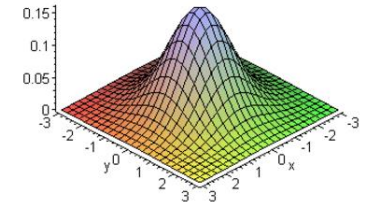
Multivariate Gaussian (Normal) Distribution

□ Multivariate Gaussian Distribution

$$p(\underline{x} | \underline{w}) = \frac{1}{Z(\underline{w})} h(\underline{x}) \exp\{\underline{w}^T \underline{T}(\underline{x})\} = h(\underline{x}) \exp\{\underline{w}^T \underline{T}(\underline{x}) - A(\underline{w})\}$$

Bivariate Normal

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



$$= \exp\left\{\text{tr}\left[-\frac{1}{2} \Sigma^{-1} (\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})^T\right] - \frac{1}{2} \ln |\Sigma| - \frac{p}{2} \ln(2\pi)\right\}$$

$$= \exp\left\{\text{tr}\left[\underline{\mu}^T \Sigma^{-1} \underline{x} - \frac{1}{2} \Sigma^{-1} \underline{x} \underline{x}^T\right] - \frac{1}{2} \Sigma^{-1} \underline{\mu} \underline{\mu}^T - \frac{1}{2} \ln |\Sigma| - \frac{p}{2} \ln(2\pi)\right\}$$

$$= \frac{1}{h(\underline{x})} \exp\left\{\text{tr}\left[\underbrace{\begin{pmatrix} \Sigma^{-1} \underline{\mu} \\ -\frac{1}{2} \Sigma^{-1} \end{pmatrix}}_{\underline{w}} \underbrace{\begin{pmatrix} \underline{x} \\ \underline{x} \underline{x}^T \end{pmatrix}}_{\underline{T}(\underline{x})}\right] - \underbrace{\left(\frac{1}{2} \underline{\mu}^T \Sigma^{-1} \underline{\mu} + \frac{1}{2} \ln |\Sigma| + \frac{p}{2} \ln(2\pi)\right)}_{A(\underline{w})}\right\}$$

□ Information form

$$p(\underline{x}; \underline{\mu}, J^{-1}) = \frac{|J|^{1/2}}{(2\pi)^{p/2}} \exp\left\{-\frac{1}{2} (\underline{x} - \underline{\mu})^T J (\underline{x} - \underline{\mu})\right\}$$

$\underline{\mu}$: mean

Σ : covariance matrix

$J = \Sigma^{-1}$: “information” matrix



Log Partition Function: MVN

$$T(\underline{x}) = \begin{bmatrix} \underline{x} \\ \underline{x}\underline{x}^T \end{bmatrix}; \underline{w} = \begin{bmatrix} \Sigma^{-1} \underline{\mu} \\ -\frac{1}{2} \Sigma^{-1} \end{bmatrix} = \begin{bmatrix} \underline{w}_1 \\ W_2 \end{bmatrix}$$

$$\begin{aligned} A(\underline{w}) &= \frac{1}{2} \underline{\mu}^T \Sigma^{-1} \underline{\mu} + \frac{1}{2} \ln |\Sigma| + \frac{p}{2} \ln(2\pi) = -\frac{1}{4} (\underline{\mu}^T \Sigma^{-1}) \left(-\frac{1}{2} \Sigma^{-1} \right)^{-1} (\Sigma^{-1} \underline{\mu}) - \frac{1}{2} \ln |\Sigma^{-1}| + \frac{p}{2} \ln(2\pi) \\ &= -\frac{1}{4} \underline{w}_1^T W_2^{-1} \underline{w}_1 - \frac{1}{2} \ln |-2W_2| + \frac{p}{2} \ln(2\pi) \end{aligned}$$

$$\nabla_{\underline{w}_1} A(\underline{w}) = -\frac{1}{2} W_2^{-1} \underline{w}_1 = -\frac{1}{2} (-2\Sigma) \Sigma^{-1} \underline{\mu} = \underline{\mu} = E(\underline{x})$$

$$\begin{aligned} \nabla_{W_2} A(\underline{w}) &= \nabla_{W_2} \left(-\frac{1}{4} \text{tr} [W_2^{-1} \underline{w}_1 \underline{w}_1^T] - \frac{1}{2} \ln |-2W_2| \right) \\ &= \nabla_{W_2} \left(-\frac{1}{4} \text{tr} [W_2^{-1} \underline{w}_1 \underline{w}_1^T] - \frac{1}{2} \ln [(-2)^p] - \frac{1}{2} \ln |W_2| \right) \\ &= \frac{1}{4} (W_2^{-1} \underline{w}_1 \underline{w}_1^T W_2^{-1}) - \frac{1}{2} (W_2)^{-1} \\ &= \frac{1}{4} ((-2\Sigma) \Sigma^{-1} \underline{\mu} \underline{\mu}^T \Sigma^{-1} (-2\Sigma)) - \frac{1}{2} (W_2)^{-1} = \underline{\mu} \underline{\mu}^T + \Sigma = E\{\underline{x}\underline{x}^T\} \end{aligned}$$

$$\nabla_{\underline{w}} (\underline{w}^T A \underline{w}) = \nabla_{\underline{w}} (\text{tr} [A \underline{w} \underline{w}^T]) = 2A \underline{w}$$

$$\nabla_A [\ln |A|] = A^{-1}$$

$$\nabla_{\underline{w}_1 \underline{w}_1}^2 A(\underline{w}) = -\frac{1}{2} W_2^{-1} = -\frac{1}{2} (-2\Sigma) = \Sigma = \text{cov}(\underline{x})$$

Rest messy!



Examples of Sufficient Statistics

Normal : $\sum_{i=1}^N x_i$; $\sum_{i=1}^N x_i^2$ for scalars; $\sum_{i=1}^N \underline{x}_i$; $\sum_{i=1}^N \underline{x}_i \underline{x}_i^T$ for vectors

(or) sample mean and sample covariance

Bernoulli: $\sum_{i=1}^N x_i$

Poisson: $\sum_{i=1}^N x_i$

Sufficient statistics help in experimenting with data transformations

Exponential : $\sum_{i=1}^N x_i$

Uniform : $\max(x_i)$

Gamma : $\sum_{i=1}^N x_i$; $\sum_{i=1}^N \ln x_i$ (or) $\sum_{i=1}^N x_i$; $\prod_{i=1}^N x_i$

Beta : $\sum_{i=1}^N \ln x_i$; $\sum_{i=1}^N \ln(1 - x_i)$ (or) $\prod_{i=1}^N x_i$; $\prod_{i=1}^N (1 - x_i)$

More at: https://en.wikipedia.org/wiki/Exponential_family



Gaussian Likelihood-Gaussian Prior for Learning the Mean

- Want to learn the mean $\underline{\mu}$ (covariance Σ known)

$$p(\underline{x}^n | \underline{\mu}) = N(\underline{\mu}, \Sigma_v) \text{ and } p(\underline{\mu}) = N(\underline{\hat{\mu}}^0, \hat{\Sigma}^0)$$

- Data: Let $D^n = \{\underline{x}^1, \underline{x}^2, \dots, \underline{x}^{n-1}, \underline{x}^n\} = \{D^{n-1}, \underline{x}^n\}$

$$p(D^n | \underline{\theta}) = p(\underline{x}^n | \underline{\theta}) p(D^{n-1} | \underline{\theta}) \dots \text{conditionally i.i.d.}$$

- Gaussian density reproduces itself \Rightarrow Posterior density is also Gaussian

$$p(\underline{x}^n) = N(\underline{\hat{\mu}}^{n-1}, \hat{\Sigma}^{n-1} + \Sigma_v)$$

Model:

$$\underline{x}^n = \underline{\mu} + \underline{v}^n;$$

$$\underline{v}^n = N(\underline{0}, \Sigma_v)$$

$$p(\underline{\mu} | D^{n-1}) = N(\underline{\hat{\mu}}^{n-1}, \hat{\Sigma}^{n-1}) \text{ and } p(\underline{x}^n | \underline{\mu}) = N(\underline{\mu}, \Sigma_v)$$

$$\underline{\hat{\mu}}^n = \underline{\hat{\mu}}^{n-1} + \hat{\Sigma}^{n-1} [\hat{\Sigma}^{n-1} + \Sigma_v]^{-1} (\underline{x}^n - \underline{\hat{\mu}}^{n-1}) \Rightarrow \text{RLS or a simple Kalman filter}$$

$$\hat{\Sigma}^n = [(\hat{\Sigma}^{n-1})^{-1} + \Sigma_v^{-1}]^{-1} = \hat{\Sigma}^{n-1} - \hat{\Sigma}^{n-1} [\hat{\Sigma}^{n-1} + \Sigma_v]^{-1} \hat{\Sigma}^{n-1}$$

$$= \hat{\Sigma}^{n-1} [\hat{\Sigma}^{n-1} + \Sigma_v]^{-1} \Sigma_v$$

Add & subtract Σ_v

$$= \Sigma_v [\hat{\Sigma}^{n-1} + \Sigma_v]^{-1} \hat{\Sigma}^{n-1}; \hat{\Sigma}^0 \text{ initial value}$$

Explicit solution:

$$\underline{\hat{\mu}}^n = \hat{\Sigma}_n \left(\hat{\Sigma}_0^{-1} \underline{\hat{\mu}}^0 + \Sigma_v^{-1} \sum_{j=1}^n \underline{x}^j \right)$$

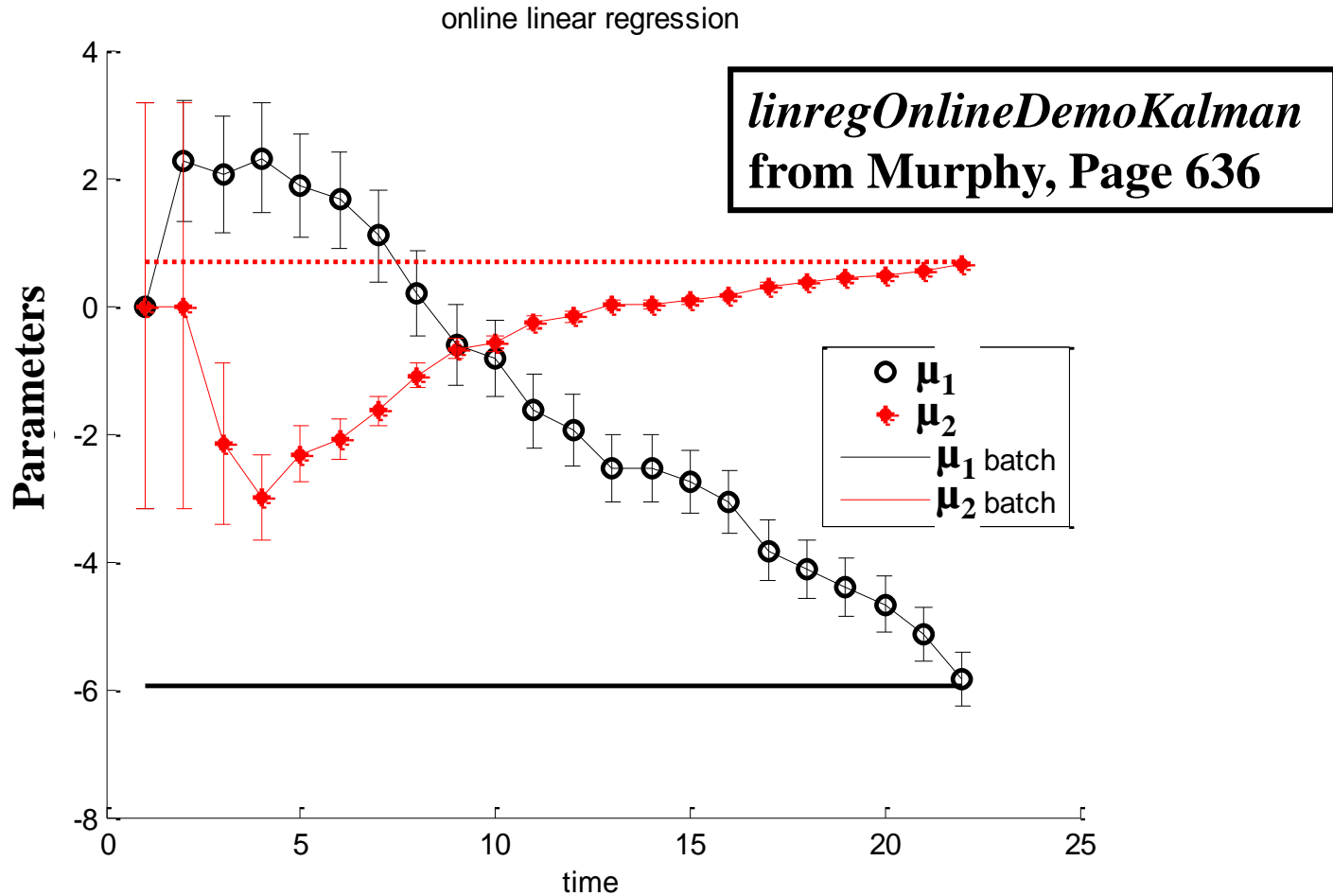
$$= \hat{\Sigma}_n \left(\hat{\Sigma}_0^{-1} \underline{\hat{\mu}}^0 + n \Sigma_v^{-1} \bar{\underline{x}}_n \right)$$

$$\hat{\Sigma}_n^{-1} = \hat{\Sigma}_0^{-1} + n \Sigma_v^{-1}$$

What if mean drifts? KF, MM, IMM, **Information filter**



Recursive Estimation of Constant Parameters



Measurement can be a linear combination of parameters corrupted by noise



Recursive Estimation of Drifting Parameters

Model :

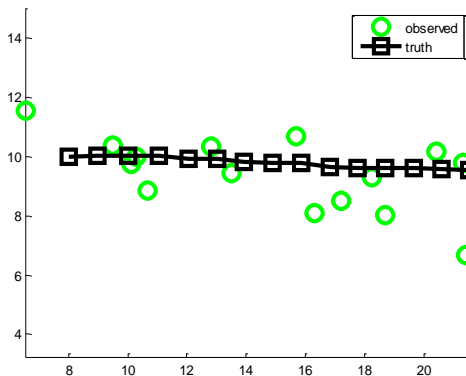
$$\underline{\mu}^{n+1} = \underline{\mu}^n + \underline{\delta}^n \Delta + \underline{w}_1^n; \Delta = \text{sampling interval}; \underline{w}_1^n \sim N(\underline{0}, Q_1); \underline{\delta}^n = \text{drift rate at sample } n$$

$$\underline{\delta}^{n+1} = \underline{\delta}^n + \underline{w}_2^n; \underline{w}_2^n \sim N(\underline{0}, Q_2)$$

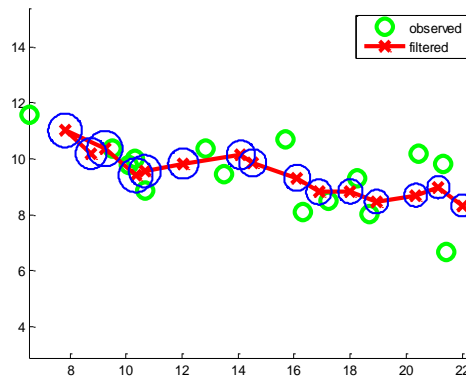
2-state nearly-constant velocity model

$$\underline{x}^n = \underline{\mu}^n + \underline{v}^n; \underline{v}^n = N(\underline{0}, \Sigma_v)$$

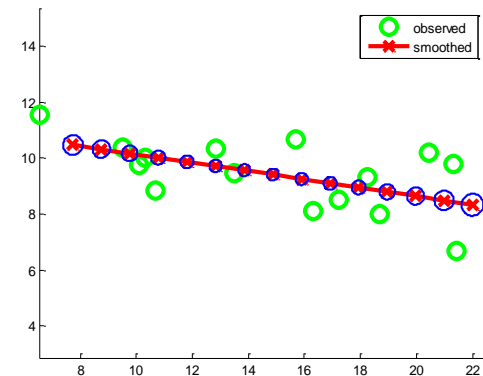
kalmanTrackingDemo from Murphy, Page 632



**True trajectory
and measurements**



Filtered Estimate



Smoothed Estimate



Entropy and Conditional Entropy

- Suppose $X \in \{x_1, x_2, \dots, x_m\}$ and $Y \in \{y_1, y_2, \dots, y_n\}$ are discrete random variables with the joint probability mass function $p(x, y)$
 - **Example:** In the text categorization problem, X is a random variable over the set of terms, and Y is a random variable over the set of documents. Evidently, $p(x, y)$ is an m by n matrix, termed a two-dimensional contingency table or a two-way frequency table; this is often estimated from the co-occurrence data.

- **Entropy of X**

$$\text{Gaussian : } H(X) = \log(\sigma\sqrt{2\pi e})$$

$$H(X) = E[-\log_2 p(X)] = - \sum_{x \in \{x_1, x_2, \dots, x_m\}} p(x) \log_2 p(x)$$

- A measure of the average amount of information (in bits) required to describe the random variable.
- **Conditional Entropy $H(X|Y)$**
 - Entropy of random variable X given another random variable, Y
- Entropy of a pair of random variables (X, Y) is the entropy of one plus the conditional entropy of the other, i.e.,

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

We find it convenient to use \ln instead of \log_2



Mutual Information (MI)

■ Mutual Information (Information Gain, MI)

- Reduction in uncertainty due to knowledge of another random variable

$$I(X;Y) = H(X) - H(X | Y) = H(Y) - H(Y | X)$$

$$= H(X) + H(Y) - H(X, Y) = I(Y; X)$$

$$= \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)} = \sum_x p(x) \sum_y p(y | x) \log_2 \frac{p(y | x)}{p(y)}$$

$$= \sum_y p(y) \sum_x p(x | y) \log_2 \frac{p(x | y)}{p(x)}$$

- X and Y are independent, i.e., $p(x, y) = p(x)p(y) \Rightarrow$ mutual information is zero
- When there is a deterministic one-to-one relationship between X and Y , mutual information is the largest (i.e., equal to the entropy of the random variable)
- In general, the mutual information is bounded by

$$0 \leq I(X;Y) \leq \min(H(X), H(Y))$$



Mutual Information Between Gaussian Variables

$$p(\underline{x}, \underline{y}) = N\left(\begin{bmatrix} \underline{\mu}_x \\ \underline{\mu}_y \end{bmatrix}; \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{xy}^T & \Sigma_{yy} \end{bmatrix}\right)$$

$$\Rightarrow p(\underline{x}) = N(\underline{\mu}_x, \Sigma_{xx}); p(\underline{y}) = N(\underline{\mu}_y, \Sigma_{yy})$$

$$p(\underline{x} | \underline{y}) = N(\underline{\mu}_x + \Sigma_{xy} \Sigma_{yy}^{-1} (\underline{y} - \underline{\mu}_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)$$

$$I(X; Y) = H(X) + H(Y) - H(X, Y)$$

$$H(X) = E_{p(\underline{x})}[-\ln p(\underline{x})] = \frac{1}{2} [n_x \ln(2\pi e) + \ln |\Sigma_{xx}|]; H(Y) = \frac{1}{2} [n_y \ln(2\pi e) + \ln |\Sigma_{yy}|]$$

$$H(X, Y) = \frac{1}{2} [(n_x + n_y) \ln(2\pi e) + \ln |\Sigma_{xx}| + \ln |\Sigma_{yy} - \Sigma_{xy}^T \Sigma_{xx}^{-1} \Sigma_{xy}|]$$

$$= \frac{1}{2} [(n_x + n_y) \ln(2\pi e) + \ln |\Sigma_{yy}| + \ln |\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T|]$$

$$I(X; Y) = \frac{1}{2} [\ln |\Sigma_{xx}| - \ln |\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T|] = \frac{1}{2} [\ln |\Sigma_{yy}| - \ln |\Sigma_{yy} - \Sigma_{xy}^T \Sigma_{xx}^{-1} \Sigma_{xy}|] \quad (1)$$

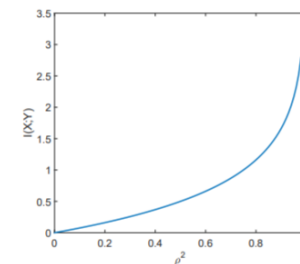
$$= -\frac{1}{2} \ln |I_{\dim(x)} - \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T| = -\frac{1}{2} \ln |I_{\dim(y)} - \Sigma_{yy}^{-1} \Sigma_{xy}^T \Sigma_{xx}^{-1} \Sigma_{xy}|$$

$$x, y \text{ scalars} : \Sigma_{xx} = \sigma_x^2; \Sigma_{xy} = \rho_{xy} \sigma_x \sigma_y; \Sigma_{yy} = \sigma_y^2$$

$$I(X; Y) = \frac{1}{2} [\ln \sigma_x^2 - \ln(\sigma_x^2 (1 - \rho_{xy}^2))] = -\frac{1}{2} \ln(1 - \rho_{xy}^2) = \ln \frac{1}{\sqrt{1 - \rho_{xy}^2}} \quad (2)$$

A greedy forward/ backward feature selection algorithm based on (1) and (2) for regression.

1. Remove redundant features (large ρ).
2. Pick features $\{x\}$ with large correlation with respect to dependent variable (output) y .





Mutual Information for Classification

$$z \in \{1, 2, \dots, C\}; \underline{x} \in R^p; p(\underline{x}) = \sum_{i=1}^C P[\underline{x}, z = i] = \sum_{i=1}^C P(z = i) p(\underline{x} | z = i) = \sum_{i=1}^C \pi_i N(\underline{x}; \underline{\mu}_i, \Sigma_i)$$

$$I(\underline{x}, z) = H(z) - H(z | \underline{x}) = H(\underline{x}) - H(\underline{x} | z)$$

It is not possible to compute $H(\underline{x})$ or $H(z | \underline{x})$ in closed-form. $H(z)$ and $H(\underline{x} | z)$ are easy to compute.

$$H(\underline{x} | z) = \frac{1}{2} \sum_{i=1}^C \pi_i (p \ln(2\pi e) + \ln |\Sigma_i|)$$

$$\text{Upper Bound on } H(\underline{x}): H(\underline{x}) = E_{p(\underline{x})} \left\{ -\ln \left(\sum_{i=1}^C \pi_i N(\underline{x}; \underline{\mu}_i, \Sigma_i) \right) \right\} \leq \sum_{i=1}^C \pi_i E_{p(\underline{x})} \left\{ -\ln \left(N(\underline{x}; \underline{\mu}_i, \Sigma_i) \right) \right\}$$

$$= \frac{1}{2} \sum_{i=1}^C \sum_{j=1}^C \pi_i \pi_j [p \ln(2\pi) + \ln |\Sigma_i| + (\underline{\mu}_j - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{\mu}_j - \underline{\mu}_i) + \text{tr}(\Sigma_i^{-1} \Sigma_j)]$$

$$\text{So, } I(\underline{x}, z) \leq \frac{1}{2} \sum_{i=1}^C \sum_{j=1}^C \pi_i \pi_j [(\underline{\mu}_j - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{\mu}_j - \underline{\mu}_i) + \text{tr}(\Sigma_i^{-1} \Sigma_j) - p] \quad (1a)$$

$$\text{When } \Sigma_i = \Sigma \text{ for all } i, I(\underline{x}, z) \leq \frac{1}{2} \sum_{i=1}^C \sum_{j=1}^C \pi_i \pi_j [(\underline{\mu}_j - \underline{\mu}_i)^T \Sigma^{-1} (\underline{\mu}_j - \underline{\mu}_i)] \quad (1b)$$

$$\text{Best single feature: } k = \arg \max_{l \in \{1, 2, \dots, p\}} \frac{1}{2} \left(\sum_{i=1}^C \sum_{j=1}^C \pi_i \pi_j \left[\frac{(\mu_{jl} - \mu_{il})^2}{\sigma_l^2} \right] \right) \quad (2a)$$

$$\text{For binary classes: } k = \arg \max_{l \in \{1, 2, \dots, p\}} \left(\frac{(\mu_{1l} - \mu_{2l})^2}{\sigma_l^2} \right) \sim d^2 \quad (2b)$$

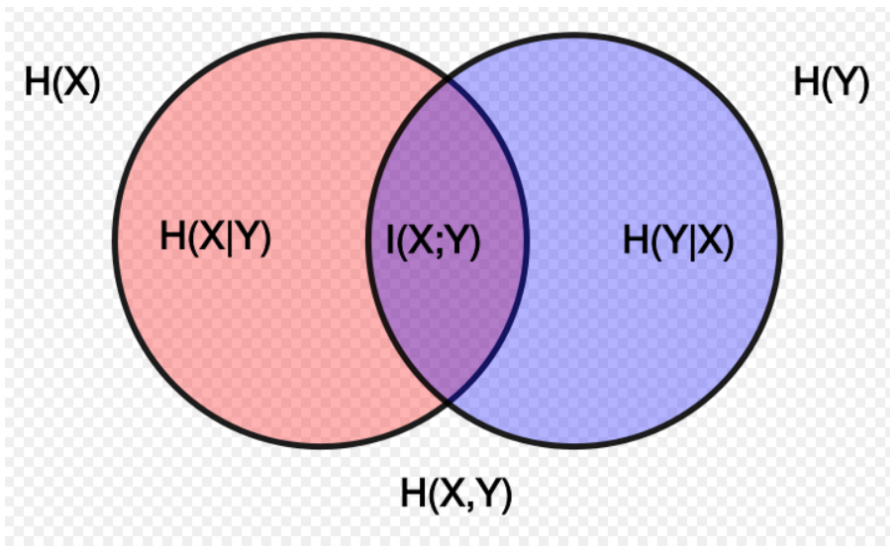
Discriminability measure

A greedy forward/backward feature selection algorithm based on (1) and (2) for classification

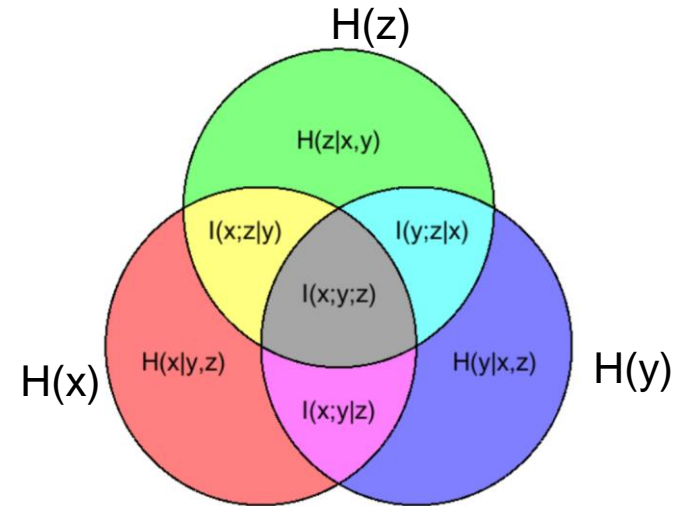


Measures of Information: Venn Diagrams

- Information (Venn) Diagram illustrates relationships among measures of information: Entropy, Joint Entropy, Conditional Entropy and MI



$$\begin{aligned}
 I(X;Y) &= H(X) - H(X|Y) = H(Y) - H(Y|X) \\
 &= H(X) + H(Y) - H(X,Y) = I(Y;X)
 \end{aligned}$$



$$\begin{aligned}
 I(X;Y;Z) &= I(X;Y) - I(X;Y|Z) \\
 &= I(Y;Z) - I(Y;Z|X) \\
 &= I(X;Z) - I(X;Z|Y)
 \end{aligned}$$

- $H(X/Y)$ is the novel information in X
- $H(Y/X)$ is the novel information in Y
- $I(X;Y)$ is the Redundant information between X and Y



Joint Mutual Information Between Gaussian Variables - 1

Consider three Gaussian random variables x, y, z

$$p(x, y, z) = N\left(\begin{bmatrix} \mu_x \\ \mu_y \\ \mu_z \end{bmatrix}; \begin{bmatrix} \sigma_x^2 & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_y^2 & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_z^2 \end{bmatrix}\right) = N\left(\begin{bmatrix} \mu_x \\ \mu_y \\ \mu_z \end{bmatrix}; \begin{bmatrix} \sigma_x^2 & \rho_{xy}\sigma_x\sigma_y & \rho_{xz}\sigma_x\sigma_z \\ \rho_{xy}\sigma_x\sigma_y & \sigma_y^2 & \rho_{yz}\sigma_y\sigma_z \\ \rho_{xz}\sigma_x\sigma_z & \rho_{yz}\sigma_y\sigma_z & \sigma_z^2 \end{bmatrix}\right)$$

$$\text{Recall } I(X;Y) = \frac{1}{2}[\ln |\Sigma_{xx}| - \ln |\Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{xy}^T|] = \frac{1}{2}[\ln |\Sigma_{yy}| - \ln |\Sigma_{yy} - \Sigma_{xy}^T\Sigma_{xx}^{-1}\Sigma_{xy}|]$$

$$x, y \text{ scalars : } \Sigma_{xx} = \sigma_x^2; \Sigma_{xy} = \rho_{xy}\sigma_x\sigma_y; \Sigma_{yy} = \sigma_y^2$$

$$\Rightarrow I(X;Y) = \frac{1}{2}[\ln \sigma_x^2 - \ln(\sigma_x^2(1 - \rho_{xy}^2))] = -\frac{1}{2}\ln(1 - \rho_{xy}^2)$$

$$\text{Recall } I(X;Y;Z) = I(X;Y) - I(X;Y|Z) = -\frac{1}{2}\ln(1 - \rho_{xy}^2) - I(X;Y|Z)$$

$$\text{Now, } I(X;Y|Z) = H(X|Z) + H(Y|Z) - H(X,Y|Z)$$

Simpler one:

$$I(X;Y|Z) = H(X|Z) - H(X|Y,Z) \dots \text{Try it!}$$

$$p(x|z) = N\left(\mu_x + \frac{\sigma_{xz}}{\sigma_z^2}(z - \mu_z); \sigma_x^2 - \frac{\sigma_{xz}^2}{\sigma_z^2}\right); p(y|z) = N\left(\mu_y + \frac{\sigma_{yz}}{\sigma_z^2}(z - \mu_z); \sigma_y^2 - \frac{\sigma_{yz}^2}{\sigma_z^2}\right)$$

$$p(x, y|z) = N\left(\begin{bmatrix} \mu_x + \frac{\sigma_{xz}}{\sigma_z^2}(z - \mu_z) \\ \mu_y + \frac{\sigma_{yz}}{\sigma_z^2}(z - \mu_z) \end{bmatrix}; \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix} - \frac{1}{\sigma_z^2} \begin{bmatrix} \sigma_{xz} \\ \sigma_{yz} \end{bmatrix} \begin{bmatrix} \sigma_{xz} & \sigma_{yz} \end{bmatrix}\right)$$

Joint Mutual Information Between Gaussian Variables - 2

$$\begin{aligned}
 I(X;Y | Z) &= \frac{1}{2} \left[\ln(2\pi e) + \ln\left(\sigma_x^2 - \frac{\sigma_{xz}^2}{\sigma_z^2}\right) + \ln(2\pi e) + \ln\left(\sigma_y^2 - \frac{\sigma_{yz}^2}{\sigma_z^2}\right) \right. \\
 &\quad \left. - 2\ln(2\pi e) - \ln \left| \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix} - \frac{1}{\sigma_z^2} \begin{bmatrix} \sigma_{xz} \\ \sigma_{yz} \end{bmatrix} \begin{bmatrix} \sigma_{xz} & \sigma_{yz} \end{bmatrix} \right| \right] \\
 &= \frac{1}{2} \ln \frac{(\sigma_x^2 - \frac{\sigma_{xz}^2}{\sigma_z^2})(\sigma_y^2 - \frac{\sigma_{yz}^2}{\sigma_z^2})}{\sigma_x^2 \sigma_y^2 - \left(\frac{\sigma_{xz}^2 \sigma_y^2 + \sigma_{yz}^2 \sigma_x^2}{\sigma_z^2} \right) + \frac{\sigma_{xz}^2 \sigma_{yz}^2}{\sigma_z^4} - \left(\sigma_{xy} - \frac{\sigma_{xz} \sigma_{yz}}{\sigma_z^2} \right)^2} \\
 &= \frac{1}{2} \ln \frac{(1 - \rho_{xz}^2)(1 - \rho_{yz}^2)}{\underbrace{1 - \rho_{xy}^2 - \rho_{yz}^2 - \rho_{xz}^2 + 2\rho_{xy}\rho_{xz}\rho_{yz}}_{\text{determinant of } 3 \times 3 \text{ correlation coeff. matrix} \geq 0}}
 \end{aligned}$$

$I(X;Y | Z) = 0$ if $\rho_{xz} = 1$
and $I(X;Y;Z) = I(X;Y)$

$$\text{So, } I(X;Y;Z) = I(X;Y) - I(X;Y | Z) = -\frac{1}{2} \ln \frac{(1 - \rho_{xy}^2)(1 - \rho_{xz}^2)(1 - \rho_{yz}^2)}{1 - \rho_{xy}^2 - \rho_{yz}^2 - \rho_{xz}^2 + 2\rho_{xy}\rho_{xz}\rho_{yz}}$$

$$\text{By symmetry, } I(X;Z | Y) = \frac{1}{2} \ln \frac{(1 - \rho_{xy}^2)(1 - \rho_{yz}^2)}{1 - \rho_{xy}^2 - \rho_{yz}^2 - \rho_{xz}^2 + 2\rho_{xy}\rho_{xz}\rho_{yz}}$$

$$I(Y;Z | X) = \frac{1}{2} \ln \frac{(1 - \rho_{xy}^2)(1 - \rho_{xz}^2)}{1 - \rho_{xy}^2 - \rho_{yz}^2 - \rho_{xz}^2 + 2\rho_{xy}\rho_{xz}\rho_{yz}}$$



Kullback-Leibler (KL) Divergence

■ Kullback-Leibler (K-L) Divergence (“Relative Entropy”)

- A measure of the distance between two pmfs, $p(x)$ and $q(x)$

$$KL(p(x) \parallel q(x)) = D(p(x) \parallel q(x)) = E_{p(x)} \left[\log_2 \frac{p(X)}{q(X)} \right] = \sum_{x \in \{x_1, x_2, \dots, x_m\}} p(x) \log_2 \frac{p(x)}{q(x)}$$

- A measure of inefficiency of assuming that the pmf is $q(x)$ when the true pmf is $p(x)$
- K-L divergence is non-negative, but not necessarily symmetric

■ Relation to Mutual Information

$$\begin{aligned} I(X, Y) &= E_{p(x,y)} \left[\log_2 \frac{p(X, Y)}{p(X)p(Y)} \right] = D(p(x, y) \parallel p(x)p(y)) = \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)} \\ &= \sum_x p(x) \sum_y p(y|x) \log_2 \frac{p(y|x)}{p(y)} = E_{p(x)} [D(p(y|x) \parallel p(y))] = E_{p(y)} [D(p(x|y) \parallel p(x))] \end{aligned}$$



K-L Divergence Between Gaussian PDFs

$$KL(p \parallel q) = -\int p(x) \ln \left\{ \frac{q(x)}{p(x)} \right\} dx = -E_p \{ \ln q(x) \} - H_p(x)$$

$$p(x) = N(\mu, \sigma^2); q(x) = N(m, s^2)$$

$$\ln p(x) = -\frac{1}{2} \ln(2\pi) - \ln \sigma - \frac{(x-\mu)^2}{2\sigma^2}$$

$$\ln q(x) = -\frac{1}{2} \ln(2\pi) - \ln s - \frac{(x-m)^2}{2s^2}$$

$$\Rightarrow KL(p \parallel q) = -E_p \left\{ \ln \frac{\sigma}{s} + \frac{(x-\mu)^2}{2\sigma^2} - \frac{(x-m)^2}{2s^2} \right\} = \ln \frac{s}{\sigma} - \frac{1}{2} + \frac{\mu^2 + \sigma^2 - 2\mu m + m^2}{2s^2}$$

$$= \ln \frac{s}{\sigma} - \frac{1}{2} + \frac{(\mu - m)^2 + \sigma^2}{2s^2} = \frac{1}{2} \left[\ln \frac{s^2}{\sigma^2} + \frac{(\mu - m)^2}{s^2} + \frac{\sigma^2}{s^2} - 1 \right]$$

$$\text{Multivariate case (} n \text{ vectors): } \frac{1}{2} \left[\ln \frac{\det(S)}{\det(\Sigma)} + (\underline{\mu} - \underline{m})^T S^{-1} (\underline{\mu} - \underline{m}) + \text{tr}(S^{-1}\Sigma) - n \right]$$

Minimum when $\underline{\mu} = \underline{m}$ and $S = \Sigma \Rightarrow$ moment matching; $p(x) = q(x)$



K-L Divergence and MI for MVN Variables

$$p(\underline{x}, \underline{y}) = N\left(\begin{bmatrix} \underline{\mu}_x \\ \underline{\mu}_y \end{bmatrix}; \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{xy}^T & \Sigma_{yy} \end{bmatrix}\right)$$

$$\Rightarrow p(\underline{x}) = N(\underline{\mu}_x, \Sigma_{xx}); p(\underline{y}) = N(\underline{\mu}_y, \Sigma_{yy})$$

$$p(\underline{x} | \underline{y}) = N(\underline{\mu}_x + \Sigma_{xy} \Sigma_{yy}^{-1} (\underline{y} - \underline{\mu}_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)$$

$$I(X; Y) = E_{p(\underline{y})} [KL(p(\underline{x} | \underline{y}) \| p(\underline{x}))]$$

$$= \frac{1}{2} E_{p(\underline{y})} \left[\ln \frac{|\Sigma_{xx}|}{|(\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)|} + \|\Sigma_{xy} \Sigma_{yy}^{-1} (\underline{y} - \underline{\mu}_y)\|_{\Sigma_{xx}^{-1}}^2 + \text{tr}(\Sigma_{xx}^{-1} (\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)) - \text{dim}(x) \right]$$

$$= \frac{1}{2} \left[\ln \frac{|\Sigma_{xx}|}{|(\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)|} + \text{tr}(\Sigma_{yy}^{-1} \Sigma_{xy}^T \Sigma_{xx}^{-1} \Sigma_{xy}) + \text{tr}(I_{\text{dim}(x)} - \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T) - \text{dim}(x) \right]$$

$$= \frac{1}{2} \ln \frac{|\Sigma_{xx}|}{|(\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)|} = -\frac{1}{2} \ln |(I_{\text{dim}(x)} - \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)| \quad \because \text{tr}(AB) = \text{tr}(BA) \ \& \ \text{tr}(I_{\text{dim}(x)}) = \text{dim}(x)$$

$$x, y \text{ scalars} : \Sigma_{xx} = \sigma_x^2; \Sigma_{xy} = \rho_{xy} \sigma_x \sigma_y; \Sigma_{yy} = \sigma_y^2$$

$$I(X; Y) = -\frac{1}{2} [\ln(1 - \rho_{xy}^2)] = \ln \frac{1}{\sqrt{1 - \rho_{xy}^2}}$$

$$I(X; Y | Z) = E_{p(z)} \{KL(p(x, y | z) \| p(x | z) p(y | z))\}$$

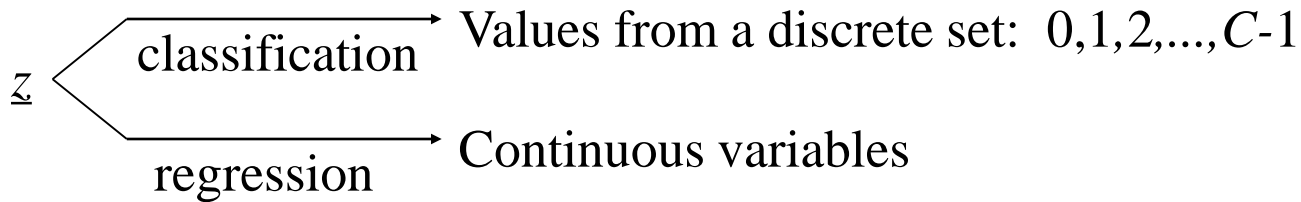
Try it for Gaussian!



Classification vs. Regression

■ Classification vs. Regression

- Set of input variables \underline{x} *observed data/features*
- Target variables \underline{z} *What you are interested in inferring*



\underline{x} can be discrete or continuous features or data

The goal is to evaluate:

- ◆ In classification $P(z = k|\underline{x}) \Rightarrow$ It is a probability (posterior probability)
- ◆ In regression $p(\underline{z}|\underline{x}) \Rightarrow$ It is a density (posterior density)



Binary Classification

□ Let us consider a binary classification case first:

$$z \in (0,1)$$

Likelihood Function

$$P(z = 1 | \underline{x}) = \frac{p(\underline{x} | z = 1)P(z = 1)}{p(\underline{x})} \quad (\text{Here, } \underline{x} \text{ is continuous.})$$

$$= \frac{p(\underline{x} | z = 1)P(z = 1)}{p(\underline{x} | z = 1)P(z = 1) + p(\underline{x} | z = 0)P(z = 0)}$$

$$= \frac{1}{1 + \frac{p(\underline{x} | z = 0) P(z = 0)}{p(\underline{x} | z = 1) P(z = 1)}}$$

Total Probability Theorem

For $a > 0, b > 0$

$$\left(\frac{a}{b}\right)^d = e^{d \ln\left(\frac{a}{b}\right)} = e^{-d \ln\left(\frac{b}{a}\right)}$$

$$= \frac{1}{1 + e^{-[\ln \frac{p(\underline{x}|z=1)}{p(\underline{x}|z=0)} + \ln \frac{P(z=1)}{P(z=0)}]}}$$

$$\ln \frac{p(\underline{x}|z=1)}{p(\underline{x}|z=0)} \quad \text{-- Likelihood Ratio}$$

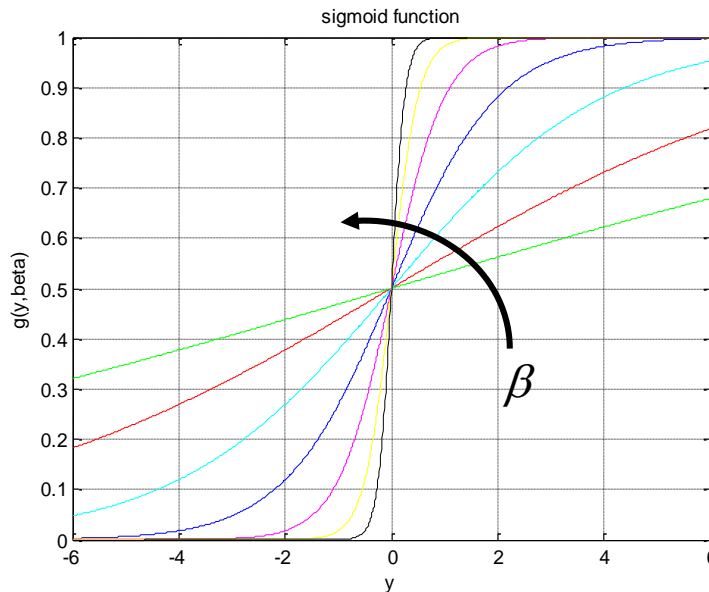
$$\ln \frac{P(z=1)}{P(z=0)} \quad \text{-- Prior Odds Ratio}$$



Sigmoid Function

$$p(z = 1 | \underline{x}) = \frac{1}{1 + e^{-y}} = g(y) \quad \text{where} \quad y = \ln \frac{p(\underline{x}|z=1)}{p(\underline{x}|z=0)} + \ln \frac{P(z=1)}{P(z=0)}$$

If $y = \underline{w}^T \underline{x}$, then this corresponds to a **simple neuron** with a **logistic sigmoid nonlinearity**. More generally,



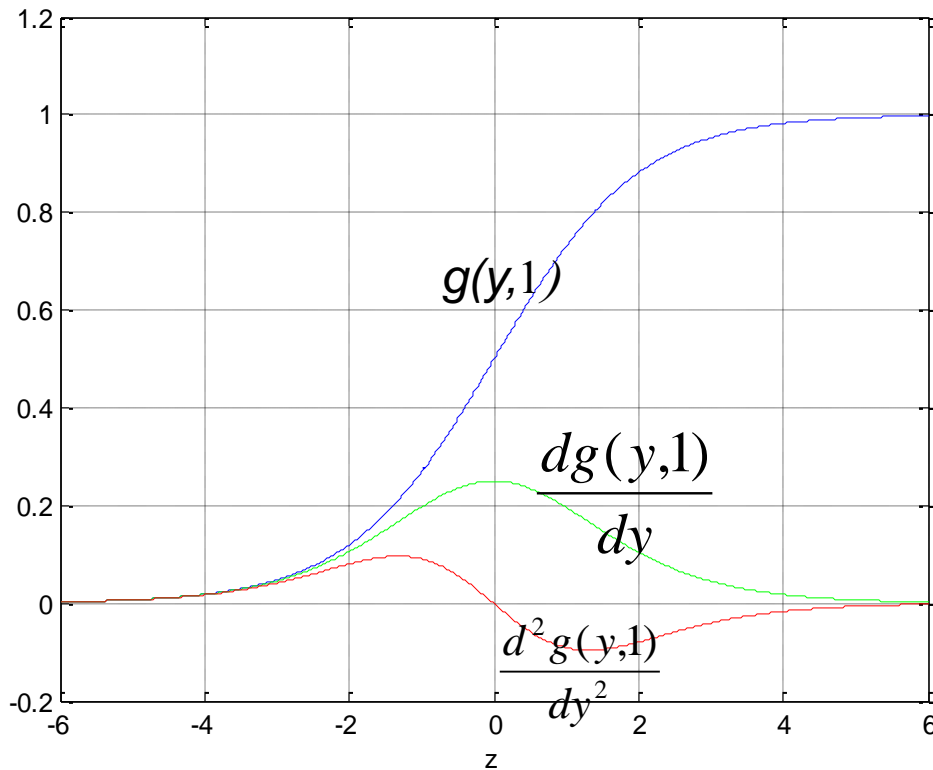
$$g(y, \beta) = \frac{1}{1 + e^{-\beta y}}$$

- What happens as $\beta \rightarrow \infty$ or $\beta \rightarrow 0$?
- Is $g(y, \beta)$ convex? **NO**. Why?

We will show that for the so-called *exponential family* of densities, this form for $y = \underline{w}^T \underline{x}$ is valid.



First and Second Derivative of Sigmoid



$$g(y, \beta) = \frac{1}{1 + e^{-\beta y}}$$
$$\frac{dg(y, \beta)}{dy} = \beta g(1 - g)$$
$$\frac{d^2g(y, \beta)}{dy^2} = \beta^2 g(1 - g)(1 - 2g)$$

$\frac{dg(y, \beta)}{dy}$ nonnegative \Rightarrow monotonic in y

$\frac{d^2g(y, \beta)}{dy^2}$ can be positive or negative \Rightarrow not convex



Discriminant Function

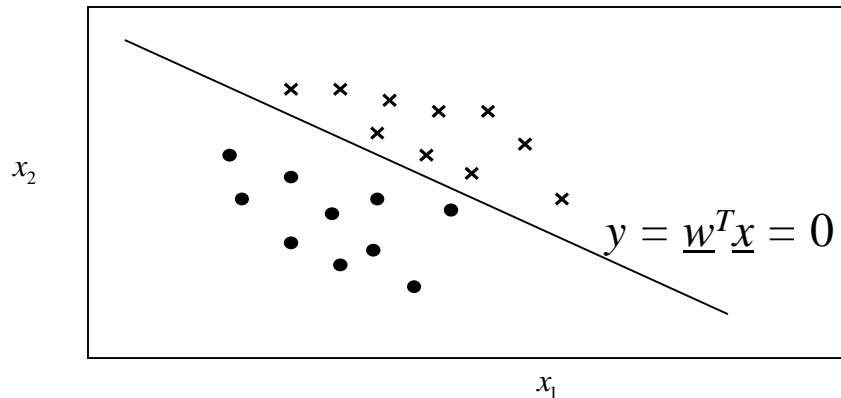
- y is also called the *discriminant function*.

$$p(z = 1 | \underline{x}) = \frac{1}{1 + e^{-y}} = g(y)$$



Any function that can be used to decide a class membership

$y = 0 \Rightarrow P(z = 1/\underline{x}) = P(z = 0/\underline{x}) \Rightarrow$ decision boundary



- Even for a non-exponential family, it is OK to use a logistic function, but they have a nonlinear form of y

$y = g(\underline{x}, \underline{w}) \Rightarrow$ Multi-layer Perceptrons (MLP)
Radial Basis Functions (RBF)

Regression

□ Regression

$p(\underline{z} | \underline{x}) \Rightarrow$ MMSE estimate is the **conditional mean**.

$$E(\underline{z} | \underline{x}) = \int_{\underline{z}} \underline{z} p(\underline{z} | \underline{x}) d\underline{z}$$

▪ For Linear Regression: $E(\underline{z} | \underline{x}) = W \underline{x} = \underline{y} = \hat{\underline{z}}$

$$\hat{\underline{z}} = \underline{y} = \sum_{j=0}^p \underline{w}_j x_j \Rightarrow \hat{z}_i = y_i = \sum_{j=0}^p w_{ij} x_j; i = 1, 2, \dots, m$$

◆ ADALINE: Adaptive Linear Element

▪ In general, we use nonlinear regression functions

$$\underline{y} = \sum_{j=0}^p \underline{w}_j \phi_j(\underline{x}) \Rightarrow y_i = \sum_{j=0}^p w_{ij} \phi_j(\underline{x}); i = 1, 2, \dots, m$$

* How to select basis functions or kernels, $\phi_j(\underline{x})$?



RBF Networks & MLP

□ RBF Networks

Radial basis function networks using Gaussian mixtures

$$\phi_j(\underline{x}) = \exp\left\{-\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1}(\underline{x} - \underline{\mu})\right\}$$

- Select $\underline{\mu}_j, \Sigma_j$
- Optimize over $\{w_{ij}\}$

⇒ Local Approximation

□ MLP: Generalize single layer perceptron

$$\begin{aligned} y_i(\underline{x}) &= g\left\{\sum_{j=1}^M w_{ij} g(\tilde{w}_j^T \underline{x})\right\} & \underline{x} &= [-1, x_1, x_2, \dots, x_p]^T \\ &= g\{w_i^T \underline{g}(\underline{x}, \tilde{w})\} \end{aligned}$$

Model of this form can approximate continuous functions to arbitrary accuracy -- universal function approximators.



Example Nonlinearities - 1

- Rectifier or rectified linear unit (ReLU). Popular in deep learning

$$g(y) = \max(0, y) = (y)^+$$

- Hyperbolic tangent

$$g(y) = \tanh(y) = (e^y - e^{-y}) / (e^y + e^{-y}) = (e^{2y} - 1) / (e^{2y} + 1) = \sigma(2y) - \sigma(-2y)$$

- Sigmoid

$$g(y) = 1 / (1 + e^{-y}) = \sigma(y)$$

- Softmax (for multiple classes as an output unit)

$$g(y_i) = e^{y_i} / \sum_{k=1}^C e^{y_k}; \sum_{i=1}^C g(y_i) = 1; g(y_i) > 0$$

- RBF unit: used with features in SVM, RBF,...

$$\phi_j(\underline{x}) = \exp\left\{-\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1}(\underline{x} - \underline{\mu})\right\}$$



Example Nonlinearities - 2

- Softplus (function approximation and Restricted Boltzmann Machine (RBM))

$$g(y) = \ln(1+e^y) = -\ln \sigma(-y) = -\ln(1 - \sigma(y))$$

soft approximation to $y^+ = \max(0, y)$

- Hard *tanh* or unit saturation block

$$g(y) = \max(-1, \min(1, y))$$

- Absolute value

$$g(y) = |y|$$

- Maxout (used with features)

$$g(\underline{x}) = \max_i (\underline{w}_i^T \underline{x} + w_{oi})$$



RBF Networks & MLP

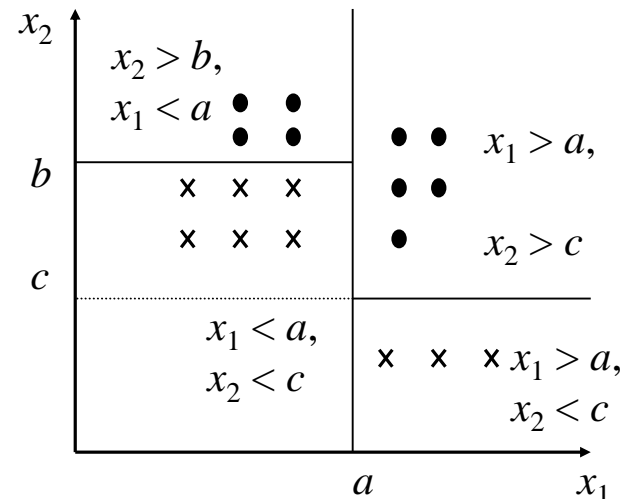
These can be extended to multiple layers $\begin{cases} \text{RBF: Mixture of exponentials} \\ \text{MLP} \end{cases}$

RBF and MLP belong to the category of *Feedforward Neural Networks* \Rightarrow Acyclic Graphs

- ⊗ Can we get insight into what is going on within MLP and RBF?
- ⊗ Can we extract rules from MLP and RBF?
- ⊗ Alternately, can we represent y_i as a function of non-overlapping basis functions



- Decision tree classifiers
- Instance-based classifiers (case-based AI reasoning)





Recurrent NNs & Density Estimation

□ Recurrent Neural Networks:

What if NNs have cycles? (\Rightarrow recurrent Neural Networks)

- Hopfield NNs
- Boltzmann and Helmholtz Machines
- Dynamic Neural Networks

□ Density Estimation:

$$p(\underline{z} | \underline{x}) = \frac{p(\underline{x}, \underline{z})}{p(\underline{x})} = \frac{p(\underline{x} | \underline{z}) p(\underline{z})}{p(\underline{x})}$$

$$p(\underline{x}) = \int_{\underline{z}} p(\underline{x}, \underline{z}) d\underline{z} \text{ for continuous } \underline{z}; \text{ (or) } \sum_{\underline{z}} p(\underline{x}, \underline{z}) \text{ for discrete } \underline{z}$$

Two key questions:

- Should we model $p(\underline{x}, \underline{z})$ or $p(\underline{z} | \underline{x})$ directly?
- What are the weights? How do you learn them?



Generative versus Discriminative

Generative

- Models $p(\underline{x}, \underline{z})$ from which $p(\underline{x})$, $p(\underline{z})$ and $p(\underline{x} | \underline{z})$ follow
- Easy to fit
- Easy to add new classes
- Handles missing data easily
- Handles unlabeled training data
- Can infer $p(\underline{x} | \underline{z})$... inverse problem
- **Hard to preprocess features**
- **Can give extreme predictions**
- **Less accurate in practice**

Discriminative

- Models $p(\underline{z} | \underline{x})$ directly ... flexibility
- **Needs optimization**
- **Retrain to add new classes**
- **Not easy to handle missing data**
- **Harder to handle unlabeled training data**
- **Not possible to infer $p(\underline{x} | \underline{z})$**
- Can preprocess features (e.g., kernels)
- Provides calibrated predictions
- More accurate in practice



Generative & Discriminative Methods

Generative

- Classifiers
 - Discriminant Analysis (DA)
 - Naïve Bayes (NB)
 - Tree-augmented NB
 - K-nearest neighbor
 - Infinite Mixture DA

Discriminative

- Classifiers (**C**) & Regression (**R**)
 - Logistic Regression (LR).. **C** only
 - Sparse Linear/Logistic Regression
 - Mixture of experts
 - MLP/RBF
 - Conditional Random Fields (CRF) ..**C**
 - Decision & Regression Trees
 - Boosting
 - Sparse kernelized linear/logistic
 - Relevance vector machine (RVM)
 - Support vector machine (SVM)
 - Gaussian processes (GP)
 - Smoothing splines ... **R** only



Popular Models of $p(\underline{x}|\underline{z})$ and $p(\underline{z})$

- ❑ $p(\underline{x})$ is a Gaussian mixture
 $p(\underline{x} | \underline{z})$ is multivariate Gaussian and $p(\underline{z})$ is discrete
- ❑ $p(\underline{x})$ is a mixture of multinomials
 $p(\underline{x} | \underline{z})$ is a product of discrete distributions and $p(\underline{z})$ is discrete
- ❑ Factor analysis
 $p(\underline{x} | \underline{z})$ is a product of Gaussians and $p(\underline{z})$ is a product of Gaussians
- ❑ Probabilistic independent component analysis (ICA)
 $p(\underline{x} | \underline{z})$ is a product of Gaussians and $p(\underline{z})$ is a product of Laplace
- ❑ Multinomial PCA
 $p(\underline{x} | \underline{z})$ is a product of discrete and $p(\underline{z})$ is a product of Gaussian
- ❑ Sigmoid belief net
 $p(\underline{x} | \underline{z})$ is a product of Bernoulli and $p(\underline{z})$ is a product of Bernoulli



MAP & ML Learning

So, among the learning processes, the difference is only in the data used for learning.

D = data

$$\mathbf{MAP:} \quad w = \arg \max_w p(w/D)$$

$$\mathbf{ML:} \quad w = \arg \max_w p(D/w)$$

$$p(w | D) = \frac{p(D | w)p(w)}{p(D)}$$

$$\begin{aligned} J(w) &= -\ln p(w/D) \\ &= -\ln p(D/w) - \ln p(w) + \underbrace{\ln p(D)} \end{aligned}$$

Independent of w

$$J(w) = -\ln p(D/w) - \ln p(w)$$

$$p(w) \text{ uniform} \Rightarrow \mathbf{MAP} = \mathbf{ML}$$



Supervised Learning

For Supervised Learning (assuming independent data sets)

$$\begin{aligned} p(D | w) &= \prod_{n=1}^N p(z_n, \underline{x}_n | w) \\ &= \prod_{n=1}^N p(z_n | \underline{x}_n, w) p(\underline{x}_n | w) \\ &= \prod_{n=1}^N p(z_n | \underline{x}_n, w) p(\underline{x}_n) \end{aligned}$$

$$-\ln[p(D | w)] = - \sum_{n=1}^N \{ \ln p(z_n | \underline{x}_n, w) + \ln p(\underline{x}_n) \}$$

$$J(w) = -\ln p(D | w) = \frac{1}{2} \sum_{n=1}^N \| z_n - g(\underline{x}_n, w) \|^2 \quad \Leftarrow \text{Gaussian Unity Covariance}$$

Thus, determining weights in supervised learning is a function minimization problem!



Classification Problem

For classification problems:

$$\begin{aligned}\ln P(z_n | \underline{x}_n, \underline{w}) &= z_n \underbrace{\ln P(z_n = 1 | \underline{x}_n, \underline{w})}_{g(y_n)} + (1 - z_n) \underbrace{\ln P(z_n = 0 | \underline{x}_n, \underline{w})}_{1-g(y_n)} \\ &= z_n \ln g(y_n) + (1 - z_n) \ln(1 - g(y_n))\end{aligned}$$

$$g(y_n) = P(z_n = 1 | \underline{x}_n, \underline{w}) = \frac{1}{1 + e^{-\beta y_n}}; y_n = \underline{w}^T \underline{x}_n$$

Cross-entropy
error function

$$\frac{\partial g(y_n)}{\partial y_n} = \frac{\partial}{\partial y_n} \left(\frac{1}{1 + e^{-\beta y_n}} \right) = \frac{\beta e^{-\beta y_n}}{(1 + e^{-\beta y_n})^2} = \beta g(y_n)(1 - g(y_n)) = \beta g_n(1 - g_n)$$

Gradient has a nice (and unified) form: for regression

$$\frac{\partial J}{\partial \underline{w}} = \nabla_{\underline{w}} J = - \sum_{n=1}^N \{z_n - g(\underline{x}_n, \underline{w})\} \nabla_{\underline{w}} g = -\beta \sum_{n=1}^N (z_n - g_n) g_n (1 - g_n) \underline{x}_n$$



Gradient for Classification

□ Classification

$$\begin{aligned}\frac{\partial J}{\partial w} &= -\sum_{n=1}^N \left[\frac{z_n}{g_n} \frac{\partial g_n}{\partial y_n} \frac{\partial y_n}{\partial w} - \frac{(1-z_n)}{(1-g_n)} \frac{\partial g_n}{\partial y_n} \frac{\partial y_n}{\partial w} \right] \\ &= -\sum_{n=1}^N \frac{z_n - g_n}{g_n (1-g_n)} \frac{\partial g_n}{\partial y_n} \frac{\partial y_n}{\partial w}\end{aligned}$$

$$\frac{\partial g_n}{\partial y_n} = \frac{\partial}{\partial y_n} \left(\frac{1}{1 + e^{-\beta y_n}} \right) = \frac{\beta e^{-\beta y_n}}{(1 + e^{-\beta y_n})^2} = \beta g_n (1 - g_n)$$

$$\text{So, } \frac{\partial J}{\partial w} = -\beta \sum_{n=1}^N \{z_n - g_n\} \frac{\partial y_n}{\partial w}$$

$$\text{If } y_n = \underline{w}^T \underline{x}_n$$

$$\frac{\partial J}{\partial w} = -\beta \sum_{n=1}^N (z_n - g_n) \underline{x}_n$$

Minor difference with regression case

⇒ Calculation of gradient is simple

⇒ Note the summation

⇒ With minor variations, leads to back propagation algorithm



Hessian

□ Classification

$$\frac{\partial J}{\partial \mathbf{w}} = \nabla_{\mathbf{w}} J = -\beta \sum_{n=1}^N \{z_n - g_n\} \underline{x}_n \Rightarrow \frac{\partial^2 J}{\partial \mathbf{w}^2} = \nabla_{\mathbf{w}}^2 J = \beta^2 \sum_{n=1}^N g_n (1 - g_n) \underline{x}_n \underline{x}_n^T$$

Cross-entropy error function is convex

□ Regression

$$\frac{\partial J}{\partial \mathbf{w}} = \nabla_{\mathbf{w}} J = -\beta \sum_{n=1}^N (z_n - g_n) g_n (1 - g_n) \underline{x}_n$$
$$\frac{\partial^2 J}{\partial \mathbf{w}^2} = \nabla_{\mathbf{w}}^2 J = \beta^2 \sum_{n=1}^N \left[(z_n - g_n) g_n (1 - g_n) (2g_n - 1) + (g_n (1 - g_n))^2 \right] \underline{x}_n \underline{x}_n^T$$

This is not necessarily convex



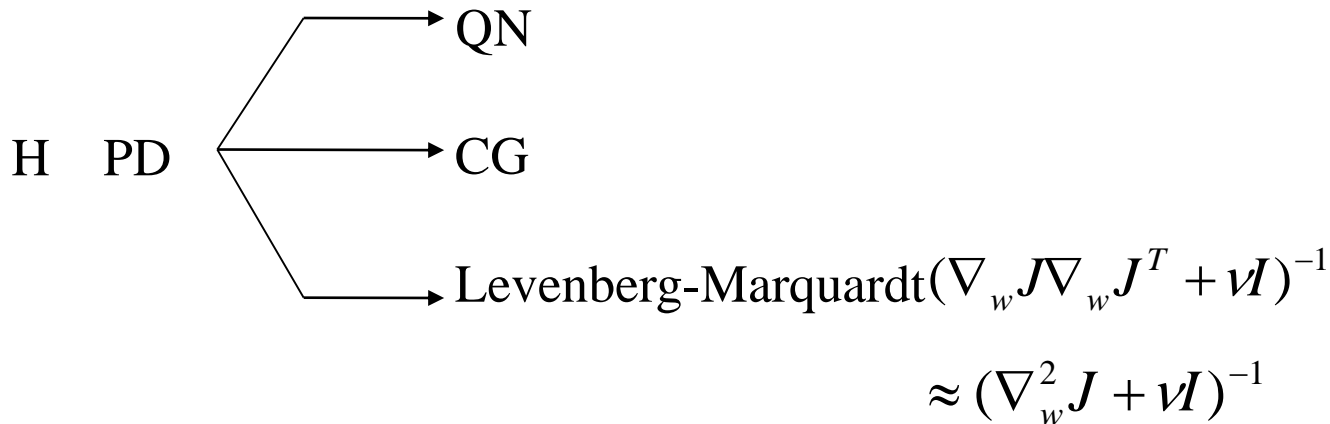
Optimization Algorithms - 1

□ Optimization Algorithms

$$\underline{w}_0 \rightarrow \underline{w}_1 \rightarrow \dots \rightarrow \underline{w}^*$$

$$\underline{w}_{k+1} = \underline{w}_k - \eta H \nabla_w J$$

$$H = \begin{cases} I & \Rightarrow \text{SD or Gradient method} \\ [\nabla_w^2 J]^{-1} & \Rightarrow \text{Newton} \end{cases}$$





Optimization Algorithms - 2

- ◆ Batch versus incremental algorithms
- ◆ Would like to estimate $Var(y_i)$
- ◆ Prune connections to simplify the architecture
- ⊗ How long to converge, how many training patterns,...

- ◆ Optimization for ill-conditioned problems (ridge regression)

Use regularization

$$\tilde{J} = J + \frac{1}{2} \nu \underline{w}^T \underline{w}$$

Regularization combats overfitting

$$\Rightarrow \nabla^2 \tilde{J} = \nabla_w^2 J + \nu I \dots \dots \text{(Levenberg-Marquardt)}$$

- ◆ L_1 -regularization (good for feature selection)

$$\tilde{J} = J + \nu \|\underline{w}\|_1$$

- Basis pursuit denoising (BPDN)
- Compressed sensing
- Lasso (Least Absolute Shrinkage & Selection Operator)



Optimization Algorithms - 3

- ◆ Combined L_1 - and L_2 -regularization

$$\tilde{J} = J + \nu_1 \|\underline{w}\|_1 + \frac{1}{2} \nu_2 \underline{w}^T \underline{w} \quad \boxed{\text{Elastic Net}}$$

- ◆ Bridge regression

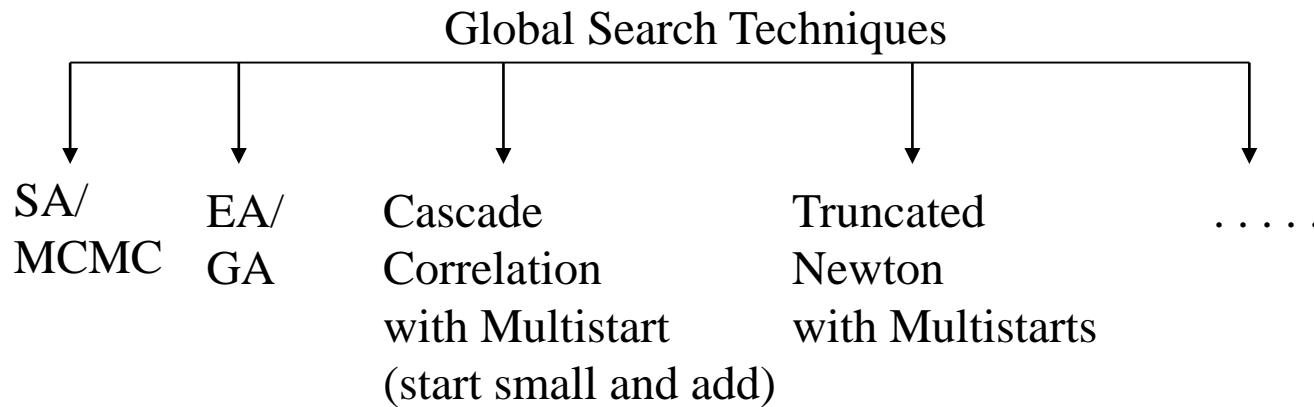
$$\tilde{J} = J + \nu_1 \sum_i |w_i|^b; b > 0$$

- ◆ Coordinate descent or probabilistic relaxation methods (EM, variational Bayes, belief propagation, empirical Bayes (Type II ML,...))



Optimization Algorithms - 4

- ◆ Numerous local minima \Rightarrow use global search techniques



Unsupervised Learning	Reinforcement Learning
Density Estimation Non-parametric techniques Kohonen SOFM, LVQ, ART Clustering (K-means, GMM, Hierarchical,...)	Application of Approximate Dynamic Programming



Summary

- ❑ Provide a systematic account of the major topics in ML based *on fundamental mathematical principles*
- ❑ Focus on pattern classification and optimization applications
- ❑ Treat classification, regression, density estimation problems in a unified way
- ❑ Generative & Discriminative classification & regression
- ❑ Learning from data: Supervised, unsupervised, semi-supervised, reinforcement, active
- ❑ Use of a variety of optimization techniques for learning weights