

02450 - F24  
Machine Learning And Data Mining  
Full Course Notes

With Prof Bjørn Sand Jensen

---

These notes thoroughly cover 100% of content found in quizzes and ~98% of total content. Please do share!

---

[Josiah Plett](#)

Immersed cause I just made a new site!

### 1 Data Feature Extraction & Visualization

Predictive → Supervised Learning  
 Descriptive → Unsupervised Learning

Classification: Discrete predictions  
 Regression: Continuous predictions

Clustering  
 Association Rule Discovery  
 Anomaly Detection

### 3 Similarity, Probabilities

$s(x, y) \rightarrow$  Similarity,  $0-1$   
 $d(x, y) \rightarrow$  Dissimilarity,  $>0$

Dissimilarity Measurements  
 $M = \text{total values}$

One-Norm  $d_1(x, y) = \sum_{j=1}^M |x_j - y_j|$   
 Euclidean  $d_2(x, y) = \sqrt{\sum_{j=1}^M (x_j - y_j)^2}$

Similarity Measurements  
 $K = \text{total values}$

SMC  $\frac{f_{00} + f_{11}}{K}$   
 Jaccard  $\frac{f_{11}}{K - f_{00}}$   
 Cosine  $\frac{x^T y}{\|x\| \|y\|}$   
 Extended Jaccard  $EJ(x, y) = \frac{x^T y}{\|x\|_2^2 + \|y\|_2^2 - x^T y}$

Transformations  
 Standardization  $\tilde{x}_{ik} = \frac{x_{ik} - \hat{\mu}_k}{\hat{\sigma}_k}$

Weighting  
 Detrending  
 Combining heterogeneous attributes

### 2 Data Features, PCA

Dataset Manipulations

Sampling: selecting subset of data  
 Feature Extraction/Transformation: New features from existing attributes  
 Feature Subset Selection:  $ABC \rightarrow AC$   
 Dimensionality Reduction: Project data down dimensions, using PCA

Common Feature Transformations

Standardization  $X_{ij} = (x_{ij} - \hat{\mu}_j) / \hat{\sigma}_j$   
 Detrending  
 Binarization  
 One-out-of-K Value  $\begin{matrix} A & B & C \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{matrix}$

Bag of words Representation  
 "a long sentence with stemming applied"  $\rightarrow$  long sentence |  
 stem |  
 unstem |  
 apply |

PCA - Principle Component Analysis  
 Projection onto lower-dimensional plane that maximizes variance

Types of Attributes

Discrete: countable set  
 Continuous: uncountable

Nominal: defined categories  
 Ordinal: rankable (eg <, >)  
 Interval: Nominal + Ordinal  
 Ratio: zero means absent

### Probabilities

Sum:  $P(A|B) + P(\bar{A}|B) = 1$   
 Product:  $P(ABC) = P(B|AC)P(A|C)$   
 Bayes:  $P(A|BC) = \frac{P(B|AC)P(A|C)}{P(B|AC)P(A|C) + P(\bar{B}|AC)P(\bar{A}|C)}$   
 Given:  $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Independence:  $P(AB) = P(A)P(B)$   
 Expectation:  $E[f] = \sum_{i=1}^N f(x_i)P(x_i)$   
 Variance:  $\text{Var}[X] = \sum_{i=1}^N (x_i - E[X])^2 P(x_i)$   
 Bernoulli:  $P(b|\theta) = \theta^b (1-\theta)^{1-b}$

Mean:  $E[X] = \sum_{i=1}^N x_i P(x_i)$

### 4 Data Visualizations + Probability Densities

Probability VS Density  
 Events are Intervals  
 Probability is the Integral  
 Density is the Slope

### Statistics

Mean:  $\bar{x} = E[X]$   
 Expected:  $E[g] = \int_{-\infty}^{\infty} g(x)P(x) dx$   
 Covariance:  $\text{cov}(x, y) = E[(x - \bar{x})(y - \bar{y})]$   
 Variance:  $\text{var}(x) = \text{cov}(x, x) = E[(x - \bar{x})^2]$   
 Standard Deviation:  $\text{std}(x) = \sqrt{\text{var}(x)}$   
 $\mu$  is mean vector,  $\Sigma$  is covariance matrix  
 Mahalanobis Distance:  $\text{mahalanobis}(x, y)^2 = (x - y)^T \Sigma^{-1} (x - y)$

### Visualizations

Single Attribute: Line plot, Histogram, Box Plot, Scatter plot, Bar plot

Relation Between Attributes: 2D Histogram, Scatterplots (Eg All Pairs of Attrs.), Matrix Plots, Parallel Coef.

### 5 Decision Trees, Regression

Unsupervised Learning

Regression: Continuous output  
 Classification: Discrete output

Data: inputs + outputs  
 Model:  $f(\text{input}) = \text{output}$   
 Cost Function:  $d(\text{output}, f(\text{input}))$

### Evaluating a Classifier

Confusion Matrix

	Actual	3	2	Accuracy $\frac{3+5}{3+5+0+2} = 0.8$
	Predicted	0	5	Error $\frac{2}{10} = 0.2$

### Evaluating a Regression Model

Average Loss per Observation  
 $E = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$   
 $L_1 = |y_i - \hat{y}_i|$   $L_2 = (y_i - \hat{y}_i)^2$

### Supervised Learning

Linear Regression

K-dimensional inputs

$f(x) = w_0 + w_1 x_1 + \dots + w_k x_k$   
 $= \sum_{i=0}^k w_i x_i = X^T W (x_0=1)$

### Classification

Decision Trees

Ask a series of questions until a conclusion is reached. Hunt's Algorithm starts at the root and repeatedly asks "good questions", making leaf nodes if all surviving data points match.

"Classification Trees", "Regression Trees"

Purity Gain  $\Delta = I(r) - \sum_{k=1}^c \frac{N(v_k)}{N(r)} I(v_k)$   
 $p(c|v) = \# \text{ class } C \in \text{branch } v$   
 Entropy  $(v) = - \sum_{c=1}^c p(c|v) \log_2 p(c|v)$   
 Gini  $(v) = 1 - \sum_{c=1}^c p(c|v)^2$   
 Class Error  $(v) = 1 - \max_c p(c|v)$

Note there are diff types of splits:  
 Binary  
 Ordinal  
 Nominal  
 Continuous

### Logistic Regression "Assume $y_i$ is binary"

$\theta_i = \sigma(f(x, w))$  where  $\sigma(z) = \frac{1}{1 + e^{-z}}$   
 $E(w) = \frac{1}{N} \sum_{i=1}^N [y_i \log(\theta_i) + (1 - y_i) \log(1 - \theta_i)]$ ,  $\theta_i = \sigma(x_i^T w)$

### General Linear Model (GLM)

$d = \text{cost function}$ ,  $g = \text{link function}$   
 $E(w) = \frac{1}{N} \sum_{i=1}^N d(y_i, g(\tilde{x}_i^T w))$

### Information Gain

$I = -\log p_i$

Residual Error  
 sum of Euclidean Loss ( $L_2$ )

### Learning

Maxima A Posteriori (MAP) Learning  
 Data:  $X = x_1 \dots x_N$   $Y = y_1 \dots y_N$   
 Parameters:  $w_i^*$  relates  $x_i \rightarrow y_i$   
 Maximizing  $w^* = \text{max}_{w,p} p(w, X, y)$  is equivalent to  
 $w^* = \text{min}_w E(w)$  where  
 $E(w) = \frac{1}{N} \sum_{i=1}^N [-\log p(y_i | x_i, w) - \log p(w)]$

MATRIX MULTIPLICATION  
 $\begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} = \begin{bmatrix} 1a + 2b + 3c \\ 2a + 3b + 4c \end{bmatrix}$

NORMS  
 Euclidean  $\|x\|_2 = \sqrt{x^T x}$   
 Frobenius  $\|A\|_F = \sqrt{\sum x_{ij}^2}$

6) Overfitting, Cross-Validation, Nearest-Neighbour

Overfitting

To solve, try to control "model complexity" eg. K in linear regression, or "pruning" in d. trees.

Errors

Training Error

$$E_{M,k}^{train} = \frac{1}{N^{train}} \sum_{i \in D^{train}} (y_i - f_{M,k}(x_i, w))^2$$

Test Error

Same as train but for test data.

Generalization Error

1. error is test error over an infinitely large test set.

If we had many test sets  $D_1, \dots, D_K$ , then:

$$E_M^{gen} \approx \frac{1}{K} \sum_{k=1}^K E_{M, D_k}^{test}$$

Bayesian Classification

EQUATIONS

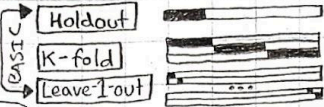
$$p(y|x_1, \dots, x_M) = \frac{p(x_1, \dots, x_M|y)p(y)}{\sum_{k=0}^{K-1} p(x_1, \dots, x_M|y=k)p(y=k)}$$

$$p(x_1, \dots, x_M|y=k) = \frac{\# \text{ of times } y=k \text{ AND we get } x_1, \dots, x_M}{\# \text{ of times } y=k}$$

Naive Bayes Assumption:  $p(x_1, x_2, \dots, x_M|y) = p(x_1|y)p(x_2|y) \dots p(x_M|y)$

Cross-Validation

Purpose: Estimate generalization Error.



1-layer Cross Validation

Purpose: Select best of S models.

Run basic cross-validation, and simply select the best one!

Forward Selection

or "Sequential Feature Selection"

1. Start w/ no features.
2. Compute cross-validation error for all +1-feature subsets.
3. Choose the best one.
4. Repeat until no change.



Backward Selection is in reverse.

2-Layer Cross-Validation

Purpose: Select optimal model AND estimate its generalization error.

Outer Loop is for comparing models via basic cross-validation

Inner Loop is for estimating model parameters via 1-layer CV

Nearest Neighbours

Can use euclidian distance or a specified distance function.

KNN with Leave-1-Out CV

For each observation  $x_i$ :

1. Temporarily remove  $x_i$ .
2. Find K nearest neighbours around  $x_i$ .
3. Determine if  $x_i$  is classified correctly based on its neighbourhood,  $c_i = 0, 1$

7) Performance Evaluation, (Naive) Bayes

Performance Evaluation

- Outline
- What's our overall Objective?
  - What tools are available?
  - What specific test should we use?

Hypothesis Testing

Let Z (perhaps  $Z = E_A^{gen} - E_B^{gen}$ ) be a quantity of interest.

Test:  $H_0: Z=0$  vs  $H_1: Z \neq 0$

1. Estimate test error somehow.
2. Calculate p-value for  $H_1$ .

Estimated result:  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n Z_i$

Statistical Tools

Too technical for me, will not be writing down.

How to Evaluate

Training on D, we obtain prediction rules  $f_{D,A}: X \rightarrow Y$  and  $f_{D,B}: X \rightarrow Y$

then we compare generalization errors:

$$E_{D,A}^{gen} = \int p(x,y) L(f_{D,A}(x), y) dx dy$$

But this is only for dataset D. To generalize, use many  $D_i$  and compute the (weighted) average. This we'll call Setup II

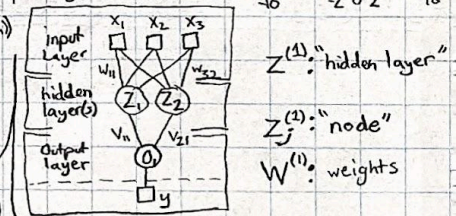
Setup II is ideal but requires lots of CV so go with Setup I if you only have a single train-test split.

Artificial Neural Networks

All same, but now model is:  $f(x) = h^{(2)}(v_0 + \sum_{j=1}^H v_j h^{(1)}(x^T w_j))$

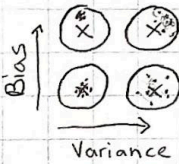
Common Choices

- $h^{(1)}(x) = \tanh(x)$
- $h^{(2)}(x) = x$
- $d(y, f(x)) = (y - f(x))^2$



8) ANNs (Artificial Neural Networks) and Bias/Variance

Bias-Variance Decomposition

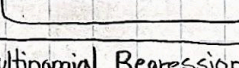


$$E_D[E^{gen}] = E_x[\text{Var}_{y|x}[y] + (\bar{y}(x) - \bar{f}(x))^2 + \text{Var}_D[f_D(x)]]$$

- Intrinsic difficulty of the problem
- Bias term; skew from true mean
- Variance term; "wigglyness"

General Linear Model

- Data  $\{X_n, y_n\}_{n=1}^N$
- Model  $f(x) = g_{link}(x^T w)$
- Cost Function  $d(y, f(x))$
- Parameters  $w = \arg \min_w \sum_{n=1}^N d(y_n, f(x_n))$



9) AUC and Ensemble Methods

- "Ensemble" → many weak classifiers combined into one strong classifier.
- "Bagging" → new datasets made by drawing randomly from pool with replacement.
- "Boosting" → train classifier on a bag, update bag weights, and repeat.

Multinomial Regression

Cost Function

(Z: is one-of-K encoding of y)

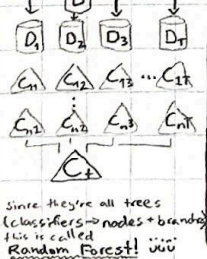
$$E = - \sum_{i=1}^N \log p(y_i | x_i)$$

Multi-class Neural Networks

suppose  $\tilde{y}_1, \dots, \tilde{y}_K$  are outputs of a NN

$$E = - \sum_{i=1}^N \log p(y_i | \tilde{y}_i)$$

Adaboost



$w_i(t) = \tilde{w}_i$  for  $i=1 \dots N$   
 for  $t=1 \dots T$  do:  
 $D_t$  by Bagging  
 $f_t$  is classifier on  $D_t$   
 $E_t = \sum_{i=1}^N w_i(t) (1 - d(f_t(x_i), y_i))$   
 $\alpha_t = \frac{1}{2} \log \frac{1 - E_t}{E_t}$   
 $w_i(t+1) = \frac{\tilde{w}_i(t+1)}{\sum_{j=1}^N \tilde{w}_j(t+1)}$ ,  $\tilde{w}(t+1) = \begin{cases} w_i(t)e^{-\alpha_t} \\ w_i(t)e^{\alpha_t} \end{cases}$

Solving the Class Imbalance Problem

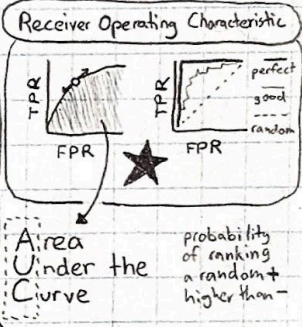
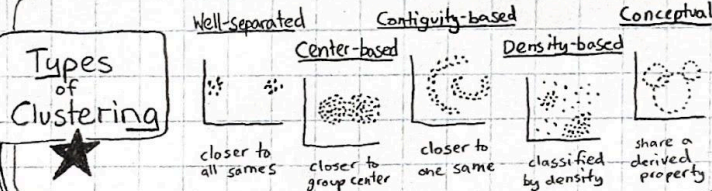
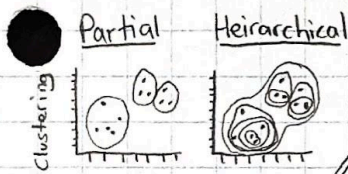
- Resampling: Undersampling, oversampling, in between.
- Modify classification algo.
- consider class imbalance when measuring performance

Precision & Recall

First, the Confusion Matrix

	Predicted		
	Positive	Negative	
Actual	Positive	FN	Precision $P = \frac{TP}{TP+FP}$
	Negative	FP	

### 10 Clustering: K-means, Hierarchical



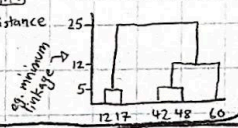
#### K-Means Clustering

- 1 Select K points as initial centroids.
- 2 Repeat:
  - Form K clusters by assigning points to closest centroids.
  - Recompute centroids.
- 3 Until: centroids don't change.

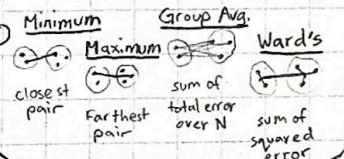
#### Agglomerative Hierarchical Clustering

- 1 Initialize proximity matrix.
- 2 Repeat:
  - Merge the closest clusters.
  - Update proximity matrix.
- 3 Until: there's just 1 cluster.

#### Dendrograms



#### Linkage Functions



### 11 Mixture Models (MM) and Density Estimation

#### Gaussian Mixture Model (GMM)

A new, probability-based method of clustering.

Locations:  $\mu_k$   
Shapes:  $\Sigma_k$   
Sizes:  $w_k$

E-step

$$p(z_n=k|x_n) = \frac{p(z=k)p(x_n|z_n=k)}{\sum_K p(z=k)p(x_n|z_n=k)}$$

M-step

$$N_k = \sum_{n=1}^N p(z_n=k|x_n) \quad w_k = \frac{N_k}{N}$$

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N x_n p(z_n=k|x_n)$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N (x_n - \mu_k)(x_n - \mu_k)^T$$

#### EM Algorithm

- 1 Select some initial models.
- 2 Repeat:
  - Expectation: calculate  $p(z_n=k|x_n)$  for all  $x_n$ .
  - Maximization: for each distribution, estimate parameters by max likelihood.
- 3 Until: params don't change.

"Mixture Models" are assessed via cross-validation, converging on a good K.

#### Probability

$\mu$  = Mean ("mu" → "M")  
 $\sigma$  = Standard Deviation ("sigma" → "S")  
 $\sigma^2$  = Variance

#### normal distribution

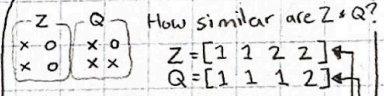
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

#### multivariate normal distribution

$$\Sigma^{-1} = \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{var}(x_2) \end{bmatrix}$$

$$N(x|\mu, \Sigma^{-1}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

#### Comparing Partitions



$$b^Z = [1 0 0 0 1]$$

$$b^Q = [1 1 0 1 0]$$

S = # of  $b^Z = b^Q$  AKA total matches  
D = # of  $b^Z \neq b^Q$  AKA total mismatches

#### Rand Index Jaccard Sim.

$$R(Z, Q) = \frac{S+D}{\frac{1}{2}N(N-1)} = \frac{3+3}{4 \times 3 / 2} = 1$$

$$J(Z, Q) = \frac{S}{\frac{1}{2}N(N-1) - D} = \frac{3}{4 \times 3 / 2 - 3} = 1$$

$$n^Z = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \quad n^Q = [3 \ 1] \quad n = \begin{bmatrix} 2 & 0 \\ 1 & 1 \end{bmatrix}$$

#### Information & Entropy Recap

Information from  $p_i$ :  $I = -\log p_i$   
Avg Information is Entropy:  $[I] = -\sum_{i=1}^N p_i \log p_i = H[P_X]$   
Mutual Info:  $MI[X;Y] = H[P_X] + H[P_Y] - H[P_{X,Y}]$   
Normalized:  $NMI[X;Y] = \frac{MI[X;Y]}{\sqrt{H[P_X]H[P_Y]}}$

### Anomaly Detection

#### Density-based Techniques

Kernel Density Estimator

$$p(x) = \sum_{n=1}^N \frac{1}{N} \mathcal{N}(x|x_n, \sigma^2 I)$$

Inverse distance density estimator

$$\text{density}_{x_i}(x_i, K) = \left( \frac{1}{K} \sum_{x \in N_{x_i}(x_i, K)} d(x_i, x) \right)^{-1}$$

Average Relative Density

$$\text{ard}(x_i, K) = \frac{\text{density}_{x_i}(x_i, K)}{\frac{1}{K} \sum \text{density}_{x_j}(x_j, K)}$$

Anomalies have very low density in multiple methods.

Inverse average distance to K nearest neighbours

ARD

### 12 Association Mining

#### Definitions

Itemset: eg. {A, C, D}, {A, D}, {}  
Support for Itemset X: % that have X  
Association Rule:  $X \rightarrow Y$  where  $X \cup Y = \emptyset$

Support: (for  $X \rightarrow Y$ )  
 $s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N} = P(X, Y)$

Confidence:  
 $c(X \rightarrow Y) = \frac{s(X \cup Y)}{s(X)} = P(Y|X)$

An algorithm for much more efficiently finding highly supported itemsets.

- 1 Find all 1-itemsets ( $L_1$ ), eg  $L_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$
- 2 Generate k-itemsets by merging rows in  $L_{k-1}$ , eg  $L_2 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$
- 3 Remove any rows in  $L_k$  that contain any subset of length  $k-1$  that is infrequent, AKA not in  $L_{k-1}$ .
- 4 Compute support of rows in  $L_k$ , remove any with insufficient support.
- 5  $L_1 \dots L_k$  are now all "frequent" itemsets.

#### A-priori Algorithm

Thank you for using my notes!  
If you want to say thanks, just share them with more people (or maybe even your prof for better reach)

*We are permitted 2 sheets of A4 paper on this exam. Use this last blank page for any personal notes as you see fit.*



*Obligatory cute red panda photo 😊*