

## Lecture 10

# Machine Learning

Let's go!!

when we want to look at specific types of errors → confusion matrix

**Accuracy:** correct preds/predictions.

**Cheat:** do well on common classes,

will automatically get high accuracy

**Misclassification rate:**  $1 - \text{accuracy}$

$0-1$  Loss matrix:

$$* \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} L(1,1) & L(1,2) \\ L(2,1) & L(2,2) \end{bmatrix}$$

Binary confusion matrix:

		pred 0
true 1	TP	FN
	FP	TN
		↑ predicted as negative predicted as positive

Possibility: collapse  $k$  classes onto 1 class

**True Positive Rate:**  $TP/P$

**False Positive Rate:**  $FP/N$

Bayes:

$$P(Y=1|x) > P(Y=2|x)$$

$$P(Y=1|x) > 1 - P(Y=2|x)$$

$$P(Y=1|x) > \frac{1}{2} \leftarrow \text{threshold}$$

Alternatives to Bayes Classifier:

Bayes: class with highest posterior prob.

Alt.: choose a different threshold

**Posterior expected loss:**  $L_k$  for class  $k$   
if we just classify EVERYTHING to class  $k$

$$L_1(x) < L_2(x)$$

$$L(2,1)P(Y=2|x) < L(1,2)P(Y=1|x)$$

$$L(2,1)(1 - P(Y=1|x)) < L(1,2)P(Y=1|x)$$

$$\frac{L(2,1)}{L(1,2) + L(2,1)} * < P(Y=1|x)$$

\* relative size of the losses  
0-1 loss has threshold 0.5

$$\frac{L(2,1)}{L(1,2) + L(2,1)} = \frac{1}{1+1} = \frac{1}{2}$$

losses

Loss matrix

$$\begin{bmatrix} L(1,1) & L(1,2) \\ L(2,1) & L(2,2) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$P_t = T_0 = \text{prob. threshold}$

$$\frac{1 - P_t}{P_t} = \frac{L(1,2)}{L(2,1)}$$

$$L(1,2) = \frac{1 - P_t}{P_t} L(2,1)$$

Precision

recall

0-1 Loss  $\rightarrow$  Bayes Classifier

$$\text{threshold} = \frac{L(2,1)}{L(1,2) + L(2,1)} = \frac{1}{1+1} = \frac{1}{2}$$

Other Losses:

threshold: between 0 and 1

matrix:  $\begin{bmatrix} 0 & 5 \\ 2 & 0 \end{bmatrix}$  ← penalizes FN more than FP

TP / predicted positives

how many predicted pos. are correct?

TP / actual positives

how many actual positives are classified correct?

precision and recall have an inherent tradeoff

F1-score

$$F_1 = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{P \cdot R}{P + R} \quad \leftarrow \text{harmonic mean of Precision and Recall}$$

harmonic mean:  $\left( \frac{1}{n} \sum_{i=1}^n x_i^{-1} \right)^{-1}$   
 $\left( \frac{a}{b} \right)^{-1} = \frac{1}{a/b} = \frac{b}{a}$   
avg. of inverse

Macro F1-score

For  $k$  classes, we have  $k$   $F_1$ -scores

macro  $F_1$ -score is the average of  $F_1$ -scores:

$$\frac{1}{K} \sum_{k=1}^K F_1\text{-score}_k$$

ROC-curve

Sensitivity against 1-specificity

TP Rate

FN Rate

this as high as possible

specific model

+ threshold varying

sensitivity

1-specificity

empty model

this as low as possible

empty model AUC =  $\frac{1}{2}$

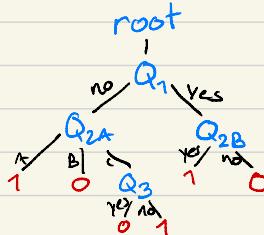
AUC: Area Under Curve

answers: which curve has least tradeoff?

empty model AUC =  $\frac{1}{2}$

# Decision Trees

we start from the root, and ask Q's to figure out which way to proceed:



**Binary Trees :** each join has exactly 2 children!

**Process :** we like shorter and more balanced trees, since performance is better!

1. Find feature with best split
2. Add to tree, and repeat on new subsets
3. Stop once criteria is reached

Falls under supervised learning

**Terminology :**

- Nodes, Edges, Root, Leaf/Terminal node - endpoints
- Internal/Interior node - Leaf & root are not these
- Decision node, Child & Parent, Descendant & Ancestors
- Subtree & Branch

- Input features can be any type
- Output can be class or numerical (Regression)
  - if numerical; entity in leaf is average

**Avoid overfitting:** end criterions! f.ex. minimum # nodes in leaf

**Assumption:** Data is somewhat clustered or otherwise grouped!

We divide featurespace into non-overlapping decision regions!  
regions are easily explainable :-)

Advantage: We only need to save the tree, and not necessarily the data.  
For balanced trees, time is  $O(\log_e(L))$  levels

Method is top-down, recursive & greedy

- starting with root node,
- splitting subsets recursively
- testing all features to decide best cut!

## Split Criterion Regression Tree:

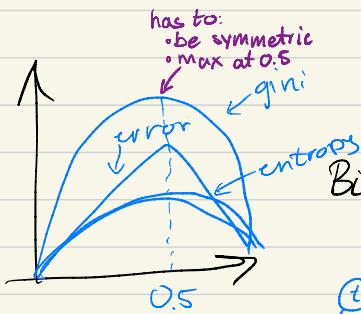
greatest possible reduction of RSS

minimize:  $R_{S1} + R_{S2}$  are splits of current subset

$$\sum_{i: x_i \in R_{S1}} (y_i - \hat{y}_{R_{S1}})^2 + \sum_{i: x_i \in R_{S2}} (y_i - \hat{y}_{R_{S2}})^2$$

↑                                   ↑

squared error Region 2  
squared error Region 1



## Binary Classification Tree

smallest possible Gini-impurity (varies 0 → 0.5)

## Impurity measures

Impurity measures  
 $\Phi(t)$ : impurity at  $t$

to be maximized  $\Rightarrow$  Goodness:  $\Phi(t) - (P_e \Phi(t_r) + P_L \Phi(t_L))$

## Gini - impurity

$$G(t) = \sum_{k=1}^K p_k(1-p_k) = 1 - \sum_{k=1}^K p_k^2$$

$$\text{Binary: } G(t) \sim 2\rho(1-\rho)$$

Minimize this

$$\text{Entropy} \quad H(t) = - \sum_{k=1}^K p_k \log p_k$$

# Decision Trees continued...

Goodness

weighted Avg.

Missing features

Stopping Criteria

Bias & Variance

Regularization  
/ shrinkage

maximize this

$$\Phi(t) - (P_R \Phi(t_R) + P_L \Phi(t_L))$$

impurity at parent minimize this 'weighted average'

$\emptyset \rightarrow$  impurity  $N \rightarrow$  # entities

$$\frac{N_L}{N_T} \emptyset_L + \frac{N_R}{N_L} \emptyset_R = P_R \Phi(t_R) + P_L \Phi(t_L)$$

right child      left child

handled by surrogate splits

↳ feature that gives most similar division to optimal

↳ we can have as many as we like, ranked

1. Stop when having fully pure leaf nodes
2. Stop when there's no more goodness
3. Stop when leaf has  $< X$  nodes
4. Stop at depth  $y$

Bias: spot on or shifted?

Variance: collected or spread out?

For decision trees:

few splits: high bias

many splits: overfit; high test error

How do we prevent overfitting?

- Early stopping/pre-pruning
  - We simply stop at a specific depth
- Pruning/post-pruning
  - We look at each possible subtree from a fully grown tree, compare performance w/o subtree and Remove subtree, if leaf has equally good perf.

# Ensemble Methods

## Wisdom of crowds

Law of Large Numbers;  $\frac{1}{m} \sum_{i=1}^m x_i \rightarrow \bar{x}$  as  $m \rightarrow \infty$

The more predictors you use, the more likely you are to arrive at correct prediction by majority vote!

Increases accuracy of model  $\hookrightarrow$

The predictors form an ensemble

BUT! predictors need to be independent

$\hookrightarrow$  in particular, errors should be independent

## Heterogenous

Different types of classifiers on same data

$\hookrightarrow$  different algorithms make diff. errors

Hard voting: Majority vote wins

Soft voting: all classifiers return  $p(k|x)$

$\hookrightarrow$  winner is class with highest avg.  $p(k|x)$

## Homogenous

Different data with same classifiers

$\hookrightarrow$  less individual point Bias/weight

Bootstrapping: Random sampling with replacement

Bagging: Bootstrap, train models in parallel, aggregate preds.

BUT! not truly independent (originate from same data)

## Out of bag eval.

because we sample with replacement, there will be datapoints that are not in our training data  $\rightarrow$  eval model on these

Bagging reduces variance

Pasting is like bagging, but w/o. replacement

# Random Forest

A forest of bagged decision trees, but!

- we only consider few features for best split (randomly chosen) INSTEAD of all ...

$M$  is the number of trees

$K$  is the number of features considered

## Extra Trees

### Extremely Randomized Trees

↳ we pick EVERYTHING at Random.

↳ at Last we compare which is best

## Boosting

Bagging does not reduce bias, so we invent boosting!

Boosting is recursively reweighting & re-training the classifier, making misclassified points have (higher weight?)  
The error of the tree is used to determine the weight of the prediction in the final vote

## Ada Boost

starting weight:  $w_i = \frac{1}{N}$  iterations:  $c=1, \dots, M$

trained models:  $h_c(x)$  for each  $c$

training error:  $E_c = \sum_{i:h_c(x_i) \neq y_i} w_i$  summing weights of wrong preds

Model weight:  $\alpha_c = \ln\left(\frac{1-E_c}{E_c}\right)$  Learning rate possible

Reweighting:  $w_i = w_i \cdot e^{\alpha_c}$  then normalize:  $w_i = \frac{w_i}{\sum_{j=1}^M w_j}$   
only misclassifications

Final Predictor:  $\text{sign}\left(\sum_{c=1}^M \alpha_c h_c(x)\right) = \text{pred}(y | x)$

## Gradient Boost

"Small steps in the right direction"

we start with a single leaf, and train a classifier on the residuals (errors) of the current model.

We use weak learners and low learning rate to avoid overfitting.

$$\text{pred} = P_{\text{init}} + LR \cdot P_{\text{res1}} + LR \cdot P_{\text{res2}}$$

## Stacking

Instead of aggregating model outputs, we train a meta-learner to blend the model outputs.  
another classifier

# Hard Margin SVM

We minimize loss by tweaking parameters in a 2D space...

i.e.  $\nabla f(x, y) = 0 \leftarrow$  slope/gradient is 0

Remember Lagrange multipliers?

Optimised when gradients of function & constraint

$$\nabla f(x, y) = \lambda g(x, y)$$

Generalisation:

if min. of  $f(x)$  within constraint  $\rightarrow$  inactive constraint

if min. of  $f(x)$  NOT within  $\rightarrow \nabla f(x) - \lambda g(x) = 0$

Lambda function  $L(x, \lambda) = f(x) - \lambda g(x)$  minimize  $x$ , maximize  $\lambda$

Karush-Kuhn-Tucker (KKT)  $g(x) = 0 \quad \lambda \geq 0 \quad g(x) \geq 0$

## Hard SVM

we only care about  
the datapoints on  
the margin...  
see slides week 15

assumption: separable binary data points

1. find widest street/band

2. middle of this is decision boundary

Method:

we look at the points and do:

$$\bar{w} \cdot x_i + b \geq 1 \quad \text{for blue points}$$

$$\bar{w} \cdot x_i + b \leq -1 \quad \text{for red points}$$

$$y_i(\bar{w} \cdot x_i + b) = 1 \quad \text{when blue=1 \& red=-1}$$

street defined by:  $y_i(\bar{w} \cdot x_i + b) - 1 = 0$

we do not know these

$$\text{width} = (x_+ - x_-) \cdot \frac{\bar{w}}{\|\bar{w}\|} = \left( \frac{1-b}{\bar{w}} - \frac{-1-b}{\bar{w}} \right) \cdot \frac{\bar{w}}{\|\bar{w}\|} = \dots = \frac{2}{\|\bar{w}\|} \leftarrow \text{maximize this width of band, by:} \\ \min_{w, b} \frac{1}{2} \|w\|^2 \quad \text{constraint } y_i(\bar{w} \cdot x_i + b) - 1 \geq 0$$

Lagrange:  $L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i^n \alpha_i [y_i(\bar{w} \cdot x_i + b) - 1]$

Lagrange continued :

$$\text{Lagrange: } L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i^n \alpha_i [y_i(w \cdot x_i + b) - 1]$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i \quad \frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$\Rightarrow L(\alpha) = \sum_i^n \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

$$\text{s.t. } \alpha \geq 0 \quad \sum_i \alpha_i y_i = 0$$

Decision Rule

$$h(x) = \sum_i \alpha_i y_i (x_i \cdot x) + b$$

# SVM Support Vector Machines

Soft margin

when we cannot separate them easily...  
But how do we decide to disregard points  
that are overlapping/outliers?

Softening constraint let some points sit between

$$y_i(w \cdot x_i + b) \geq 1$$

Slack variable

a variable, each point has  $\xi_i$ , now constraint is:

$$y_i(w \cdot x_i + b) \geq 1 - \xi_i$$

$\xi_i > 0$  for points between margins,  $= 0$  for correct

Penalize slack variable:

$$\min_{w,b} \underbrace{\frac{1}{2} \|w\|^2}_{\text{max street width}} + C \sum_{i=1}^n \xi_i \quad \text{const. } y_i(w \cdot x_i + b) \geq 1 - \xi_i \\ \text{minimize slack} \quad \xi_i \geq 0$$

Now we also care about the things IN the street

Lagrange Result:

$$\max_{\alpha} \sum_i^n \alpha_i - \frac{1}{2} \sum_i^n \sum_j^n \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

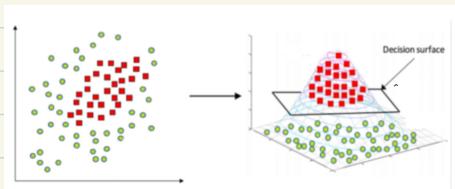
w. constraint:  $0 \leq \alpha \leq C$  &  $\sum_i^n \alpha_i y_i = 0$

Decision Rule

$$h(x) = \sum_{i \in S} \alpha_i y_i (x_i \cdot x) + b \quad b = \frac{1}{N_S} \sum_{i \in S} (y_i - \sum_{j \in S} \alpha_j y_j (x_i \cdot x_j))$$

C constant defining how much we care  
about error... hyperparameter tuning on this

Non-Linear boundaries if we can map to a higher dimension?



we define a transformation  $\phi(x)$   
↳ becomes slow fast!

Kernel Function:

$O(d)$  instead of  $O(2^d)$

$$\phi(x) \circ \phi(z) = \prod_{k=1}^d (1 + x_k z_k)$$

For training

$$\begin{aligned} \max_{\alpha} \quad & \sum_i^n \alpha_i - \frac{1}{2} \sum_i^n \sum_j^n \alpha_i \alpha_j y_i y_j (\phi(x_i) \cdot \phi(x_j)) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \quad \text{and} \quad \sum_i^n \alpha_i y_i = 0 \end{aligned}$$

$$\begin{aligned} \text{For training} \\ \arg \max_{\alpha} \quad & \sum_i^n \alpha_i - \frac{1}{2} \sum_i^n \sum_j^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \quad \text{and} \quad \sum_i^n \alpha_i y_i = 0 \end{aligned}$$

Decision boundary (for predictions):

$$h(x) = \sum_i^n \alpha_i y_i (\phi(x_i) \cdot \phi(x)) + b$$

Decision boundary (for prediction):

$$h(x) = \sum_i^n \alpha_i y_i K(x_i, x) + b$$

$$b = \frac{1}{N_S} \sum_{i \in S} (y_i - \sum_{j \in S} \alpha_j y_j K(x_j, x))$$

Kernels

Kernels is a dot-product between 2 points after they have been transformed

# Gradient Descent

Main Idea

Batch GD

Find Gradients for all datapoints and average/sum them... This is now the Gradient to follow for GD

Stochastic GD

Find Gradient for a single random datapoint, and use this for GD

Mini-Batch GD

Randomly choose a subset of the dataset, and do same as batch.

# Basics Artificial Neural Networks

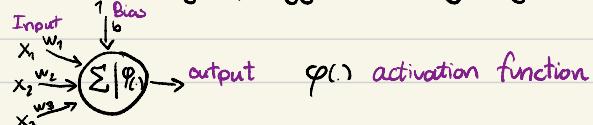
## Perceptron

### I/O

### $\varphi(\cdot)$ Activation

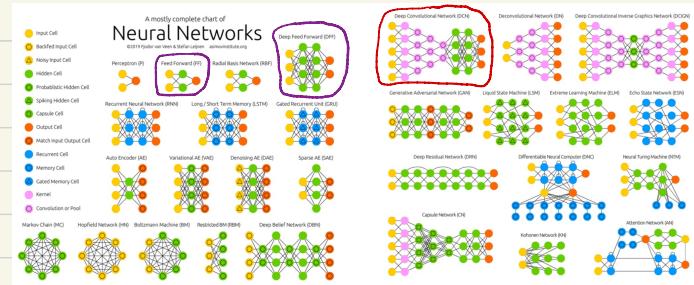
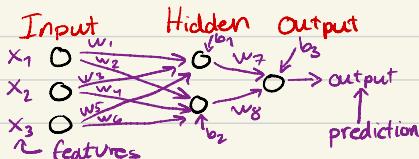
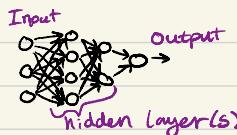
## FF Feed Forward

Inspired by biological neural network  
 ↳ receives a signal, triggers sending signal(s)



Binary: 0 or 1 ← is it active or not?

instead of 0-1 I/O, ↳ sigmoid, tanh, ReLU  
 we can do differentiable



## Output $\varphi(\cdot)$

Regression : Linear function

Binary : Sigmoid

Multiclass : Softmax

$$P(Y=1|x) \rightarrow P(Y=1|x) \rightarrow \text{after normalizing}$$

Gradient Descent  
for weights

$w_1, \dots, w_n$  &  $b_1, \dots, b_m$  are parameters.

Randomly init of parameters, GD until convergence

# Forward & Back Propagation

(n) hidden  
(o) output

**Forward pass** (starting from input layer toward output layer): compute the output and then the loss function on the training samples.

**Back propagation** (starting from output layer toward input layer): use the computed loss for updating the values of W and b.

Data is passed on, when received:

- ◆ Summed  $Z_n^{(n)} = \sum x_i w_{i,n}^{(n)} + b_n$
- ◆ activated  $a_n^{(n)} = f(Z_n^{(n)})$
- ◆ passed on to next layer

**Vectorized (matrix form):**

$$\text{Input layer } x_{1 \times 3} = [x_1, x_2, x_3]$$

$$\text{Hidden layer } Z_{1 \times 4}^{(h)} = [z_1^{(h)}, z_2^{(h)}, z_3^{(h)}, z_4^{(h)}] = x_{1 \times 3} w_{3 \times 4}^{(h)} + b_{1 \times 4}^{(h)} \quad a^{(h)} = f(Z^{(h)})_{1 \times 4}$$

$$\text{Output layer } Z_{1 \times 1}^{(o)} = [z_1^{(o)}] = a_{1 \times 4}^{(h)} w_{4 \times 1}^{(o)} + b_{1 \times 1}^{(o)} \quad a^{(o)} = g(Z^{(o)})_{1 \times 1}$$

Automatic Different.

Back-propagation is the automated diff. in reverse mode; it is RECURSIVE!

$$\frac{\partial Y}{\partial x} = \frac{\partial f}{\partial g(x)} \cdot \frac{\partial g}{\partial x}$$

inner diff. ← what to find now  
outer diff. ← previous step in recursive

### III

# Artificial Neural Networks

What if we get stuck in ?

Momentum Gradient;  
Factor in previous gradients. Becomes a weighted average, with higher weight the more recent the gradients.

Nesterov Momentum optimization

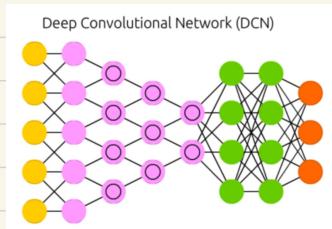
$$\begin{aligned} v &:= \beta v + \alpha \nabla L(w - \beta v) \\ w &:= w - v \end{aligned}$$

Prior gradients                                  current gradient loss

Convolutional NN

So far; layers were fully connected  
Now; Transform large input into  
smaller set of inputs to Neural Network!

Filtering



# Dimensionality Reduction

How to "throw something away"  
Drawing 3D in 2D is actually a projection,  
which offers us a reduction of dimensions!

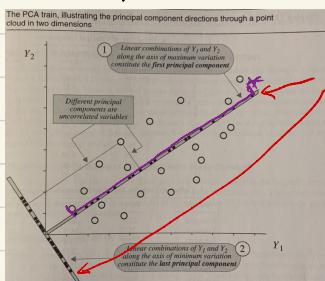
## PCA

### Principle Component Analysis

1. To represent data in a useful way
2. See if data dimensions can be reduced

No class Labels!

We prefer features with high variance



1 axis for each feature  
axes are Principle Components  
Starting with the direction of the highest variance\*

\* Scores  
 $\|a_1\| = 1$

$$A^T X = (X^T A)^T$$

The transformed values in new dimensions

$$a_1^T X = a_{11} X_1 + a_{12} X_2 + \dots + a_{1p} X_p$$

$a_1$  first column vector of  $A$

$X$ 's columns are feature vectors

Find max Var

Lagrange

$$\text{We maximise } \text{Var}(a_1^T X) = a_1^T S a_1 \quad \text{sample variance}$$

under constraint  $a_1^T a_1 = 1^*$

$$F(a_1) = a_1^T S a_1 - \lambda_1 (a_1^T a_1 - 1)$$

$$\text{Partial Derivative: } \frac{\partial F}{\partial a_1} = 2 S a_1 - 2 \lambda_1 a_1 = 2(S - \lambda_1 I) a_1$$

$$(S - \lambda_1 I) a_1 = 0 \Rightarrow S a_1 = \lambda_1 a_1$$

$$\hookrightarrow \text{Var}(a_1^T X) = a_1^T S a_1 = a_1^T a_1 = 1$$

$a_1$  is an eigenvector of  $S$  with eigenvalue  $\lambda_1$

Solution; We choose  $a_1$  with largest  $\lambda_1$

## Lagrange 2

maximise  $\text{Var}(a_2^T X) = a_2^T S a_2$

under constraint  $a_2^T a_2 = 1$  and  $a_1^T a_2 = 0$   
non-scaled features perpendicular

$$F(a_2) = a_2^T S a_2 - \lambda_1 (a_2^T a_2 - 1) - \mu (a_1^T a_2)$$

$$\frac{\partial F}{\partial a_2} = 2S a_2 - 2\lambda_1 a_2 - \mu a_1$$

Set the partial derivatives to 0

$$0 = 2S a_2 - 2\lambda_2 a_2 - \mu a_1$$

Left-multiply this by  $a_1^T$ :

$$\begin{aligned} 0 &= 2a_1^T S a_2 - 2\lambda_2 a_1^T a_2 - \mu a_1^T a_1 \\ &= a_1^T \lambda_1 a_2 - 0 - \mu \cdot 1 \\ &= \lambda_1 a_1^T a_2 - \mu \\ &= \lambda_1 0 - \mu \quad \Rightarrow \mu = 0 \end{aligned}$$

Since  $\mu = 0$ , we just need the next biggest eigenvalue as  $\lambda_2$

By induction we can prove that we just need to use the eigenvalues in desc order  $\heartsuit$

eigendecomposition  
diagonalisation

$\Delta$  is D with  $\lambda_1 \geq \dots \geq \lambda_p \geq 0$

$$\text{now: } S = A \Delta A^T$$

and now:  $A = [a_1, a_2, \dots, a_p]$  is P

Total Variance

Sum of variances  $\sum_{i=1}^p \lambda_i$

Explained Variance

For k<sub>th</sub> principle component:

$$\frac{\lambda_k}{\sum_{i=1}^p \lambda_i}$$

Variance/explained variance can be plotted (cum. we can keep m% of var)