

## 1. Basics

Observation Matrix:  $\mathbf{X} \subset \mathbb{R}^{p \times n}$  with  $p$ -dimensional random variable and  $n$  observations ( $i$ : variable,  $j$ : observation)  
 Expected Value  $\mu_X = \sum_{i \in \mathcal{X}} x_i \cdot P_X(x_i) = E[X] \in \mathbb{R}^p$   
 Variance  $\Sigma = \text{Var}(X) = E[(X - \mu)(X - \mu)^\top] \in \mathbb{R}^{p \times p}$   
 $\Sigma$  is symmetric ( $\Sigma^\top = \Sigma$ ) and positive semidefinite:  $\forall x: x^\top \Sigma x \geq 0$  or  $\Sigma \geq 0$   
 Estimated Value based on mean of observations  $n$ :  
 $\hat{\mu}_p = \frac{1}{n} \sum_{j=1}^n x_{pj}$  (calculate for every dimension  $p$ )  
 Centered Observation Matrix  $\hat{X} = X - \hat{\mu} \otimes [1 \dots 1]$   
 Covariance Matrix  $\hat{\Sigma} = \frac{1}{n-1} \hat{X} \hat{X}^\top \approx \frac{1}{n} \hat{X} \hat{X}^\top$   
 Transpose Matrix  $(XY)^\top = Y^\top X^\top$

### 1.1. Random Variables

Probability:  $Pr(X \in \mathcal{X}) = \int_{\mathcal{X}} p_X(x) dx = \sum_{\{i|x_i \in \mathcal{X}\}} p_i$

Marginal Density:  $p_X(x) = \int_{\mathbb{R}^k} p_{X,Y}(x, y) dy$

Conditional Density:  $p_{Y|X=x}(y) = \frac{p_{X,Y}(x, y)}{p_X(x)}$ , with  $\sum_y p_{Y|X=x}(y) = 1$

Expectation Value:  $E[X] = \int_{\mathbb{R}^p} x p_X(x) dx = \mu_X$

(Co)Variance:  $\text{Var}[X] = E[(X - \mu_X)(X - \mu_X)^\top]$

## 2. Statistical Decision Making

### 2.1. Loss Function

- Quadratic Loss Function:  $L(Y, f(X)) = (Y - f(X))^2$
- Minimize Expected Prediction Error:  $EPE(f) = E[L(Y, f(X))]$
- If using quadratic loss function, conditional mean is best
- If using absolute loss  $L(Y, f(X)) = |Y - f(X)|$ , conditional median is best

### 2.2. Decision Making

Map an input vector to an output value and find a decision function that minimizes the deviation between a target and the output, i.e. the loss.

- Expected Prediction Error** as a function of  $f$  is formulated by the expectation of the loss function, whose values change depending on the decision function  $f$ :  $EPE(f) = E[L(Y, f(X))]$ .
- The aim of **global methods** is to find the best explicit global function  $\hat{f} = \argmin_{f \in \mathcal{F}} EPE(f)$ .
- The aim of **local methods** is to find the best output  $\hat{c} = \argmin_{c \in \mathbb{R}} E_{Y|X=x} L(Y, c)$ , i.e. minimizes the expectation of loss on the distribution of  $Y$  conditioned on known  $X$  (only samples in region of interest).
- Binary Decision Making** can be visualised via a Joint PDF of  $X$  and  $Y$ , where  $Y$  takes on -1 or 1 and  $X$  is a continuous random variable whose distribution changes conditioned on  $Y$ .
- We can obtain the Expected Prediction Error and Expected Loss due to the assumption that we *know the Joint PDF* of  $X$  and  $Y$ , i.e. the stochastic behaviour of the system.
- For local methods with loss specified as the squared loss, expected loss is found to be the conditional expectation of  $Y$  on the distribution of  $Y$  conditioned on  $X = x$ :  $f(X) = E_{Y|X=x}[Y]$ .
- The problem with local methods is most of the time  $X$  is a continuous random variable and thus values here  $X = x$  are impossible, thus we take the **k-nearest neighbours** to  $x$  as approximations of  $x$ . This gives us a set of samples representing the distribution of values  $X = x$  such that local methods for finding the output can now be applied to these samples.

### 2.3. Generalization Error

Difference between Empirical Loss and Expected Loss for a function  $f$ . In practice, the difference between the training and the validation error.

- $G_N(f) = E[L(Y, f(X))] - \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$
- $G_N \rightarrow 0$  for  $N \rightarrow \infty$

### 2.4. k-nearest Neighbors

Problem: among many observations  $(x_i, y_i)$  there is probably none with  $x_i = x$ .

- Solvable by taking all observations into account which are near the desired  $x$ :  $\hat{f} = \text{average}(y_i | x_i \in N_k(x))$
- If the set of neighbors  $k$  increases, the average tends towards the expectation value
- $x_i$  comes closer to  $x$  if number of observations  $N$  increases
- In practice  $N$  is limited, thus imposing a model assumption for linear regression or "reducing" the dimension of  $X$  is needed.

### 2.5. Curse of Dimensionality

Most of the time it is desired that the dimensionality of samples are decreased due to various problems involved with working in high dimensions. With increasing dimensions  $p$ :

- Noise increases (accumulates over dimensions)
- The number of samples required to obtain an accurate representation of the probability distribution also increases a lot/exponentially.
- Empty Space Phenomenon**: High dimensional spaces are sparse and samples tend to be located at the tail end of their distributions.
- Distance between a point and any possible samples becomes more and more equidistant.
- For a random vector  $X \in \mathbb{R}^p$ , the probability of at least one  $X_i$  being further than  $\beta$  away from the center:

$$Pr(\|X\|_2^2 \geq \beta) = 1 - Pr(X_1^2 < \beta)P.$$

For large  $p$  this becomes 1 and it's difficult to estimate the underlying distribution.

### 2.6. Data Preparation

- Nominal Categories - No ordering, Ordinal Categories**: Ordering
- Num to Cat: Discretization, Cat to Num: Binarization
- Bag of Words**: Frequent and Distinct = strong Weight
- Frequent: Term Frequency
- Distinct: Inverse Document Frequency
- Text Prep: Remove HTML, lower case, Remove punctuation/numbers/common words, split into words

## 3. Convex Functions

**Definition Convexity**:  $C \subset \mathbb{R}^n$  a convex set, i.e for any pair  $\mathbf{x}_1, \mathbf{x}_2 \in C$  the point  $t\mathbf{x}_2 + (1-t)\mathbf{x}_1 \in C \forall t \in [0, 1]$ . A function is called **convex** if  $\forall \mathbf{x}_1, \mathbf{x}_2 \in C, t \in [0, 1]$ :

$$tf(\mathbf{x}_2) + (1-t)f(\mathbf{x}_1) \geq f(t\mathbf{x}_2 + (1-t)\mathbf{x}_1)$$

- By definition, a function is convex if a line can be drawn between any two points in a function, and all vertical values of the line are greater than or equal to all vertical values of the function between the same two points.
- A function is also convex if its second derivative is positive, or equivalently for vector inputs its Hessian matrix contains all non-negative eigenvalues.
- Several operations between convex functions also result in another convex function, including the max between two convex functions:

if  $f$  and  $g$  are convex, then  $\max(f, g)$ ,  $f + g$  and  $g \circ h$  (if  $g$  non-decreasing) are convex. A local minimum of a *strictly* convex function is a global minimum and unique.

## 4. Logistic Regression

### 4.1. Formulation via Optimization

- The aim of logistic regression to perform binary classification on a random vector, that is, map a random vector to either -1 or 1.
- In order to ensure that the outputs in logistic regression are either -1 or 1, the output in logistic regression is taken to be the sign of the output from a prediction function  $f$  which contributes to the decision.
- The loss function in logistic regression is taken to be the number of misclassifications of a prediction, that is, 1 if the decision from logistic regression does not match the target, and 0 if it does match.

$$f(z) = \begin{cases} 1 & \text{if } Y \cdot \text{sign}(f(x)) \leq 0 \\ 0 & \text{otherwise} \end{cases}$$

- It is clear that we must find an appropriate  $f$  to complete the decision function, and thus we define the global method with the **misclassification loss** as the problem we need to solve.  $\frac{1}{n} \sum_{i=1}^n L_{0,1}(y_i, f(x_i))$
- The misclassification loss is extremely difficult to solve due to being non-continuous and non-convex, and so first we need to approximate this function as a convex function since they are easy to optimize.
- A convex approximation of misclassification loss is  $\log(1 + \exp(-t))$ , which when concatenated with an affine function  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$  is also convex, and can therefore be easily optimized via minimization.

$$\min_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i(\mathbf{w}^\top \mathbf{x}_i + b)))$$

- Noting that the loss is in **negative log-likelihood** form, we can convert this to a probability by taking the exponential of likelihood.

$$Pr(Y = y|x) = \frac{1}{1 + \exp(-y(\mathbf{w}^\top \mathbf{x} + b))}$$

- Gradient-based methods can be used to find affine weights  $\mathbf{w}$  and bias  $b$  that solve the optimization problem (*first box*).
- Classify new sample:  $\text{sign}(\mathbf{w}^\top \mathbf{x}_{new} + b)$ .
- Decision boundary of LR is always a hyperplane (i.e. an affine space) in  $\mathbb{R}^p$ .
- Retraining with a new training sample affects the decision boundary.

### 4.2. Overfitting on Logistic Regression

- When the samples are linearly separable (*a dividing hyperplane exists*), there is a constraint that the target times the affine output is always greater than zero Under this constraint, no global minimum exists.

If the data is linearly separable, we can find  $(\mathbf{w}_s, b_s)$  so that:  
 $\forall i: y_i(\mathbf{w}_s^\top \mathbf{x}_i + b_s) < 0$   
 then the loss-function has no global minimum in  $\mathbb{R}^{p+1}$

This can be fixed by introducing a regulator which penalizes the magnitude of  $(\mathbf{w}, b)$ , e.g. by fixing a real constant  $\lambda > 0$  and adjusting the original cost function as  $\tilde{F}(\mathbf{w}, b) = F(\mathbf{w}, b) + \lambda \|\mathbf{w}\|^2 + b^2$ .

### 4.3. Alternative Approach (Statistics)

Linear Model:  $a = w_0 + w_1 x_1 + \dots + w_n x_n = \mathbf{w}^\top \mathbf{x}$

Probability:  $\sigma(a) = \frac{1}{1 + e^{-a}}$   $D = \{(\mathbf{x}_i, z_i)\}$

Find  $\mathbf{w}_{MLE} = \arg \max_{\mathbf{w}} P(D|\mathbf{w})$   $\mathbf{x}_i \in \mathbb{R}^d, z_i \in \{0, 1\}$

$$P(D|\mathbf{w}) = \prod_{i=1}^n \sigma(\mathbf{w}^\top \mathbf{x}_i)^{z_i} (1 - \sigma(\mathbf{w}^\top \mathbf{x}_i))^{1-z_i}$$

$$L(\mathbf{w}) = -\log(P(D|\mathbf{w})), \nabla_{\mathbf{w}} L(\mathbf{w}) = \mathbf{X}(\sigma(\mathbf{X}^\top \mathbf{w}) - \mathbf{z})$$

(Requires Hessian Matrix, which is positive semi-definite)

Results are identical up to the factor  $1/n$

## 5. Kernel Method

In its simplest form, the kernel trick means transforming data into another dimension that has a clear dividing margin between classes of data. The Kernel trick simply computes the inner products between the images of all pairs of data in the feature space to learn a nonlinear function or decision boundary.

A positive semidefinite Kernel is a function  $\kappa: \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$  such that for all finite sets  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  the *Gram-Matrix*  $\mathbf{K}_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$  is symmetric and positive semidefinite

- Symmetric  $\forall i, j: \kappa(\mathbf{x}_j, \mathbf{x}_i) = \kappa(\mathbf{x}_i, \mathbf{x}_j)$
- Positive semidefinite  $\forall \mathbf{x}: \mathbf{x}^\top \mathbf{K} \mathbf{x} \geq 0$  or all eigenvalues  $\geq 0$
- We can determine whether a particular function is a kernel by testing for violations of symmetry or positive semidefiniteness. Symmetry can easily be tested by substituting named variables into the kernel, then flipping around the variables and seeing whether the kernel expression is the same.
- Positive semidefiniteness can be tested as violated by definition, or if any diagonal values are ever negative. Furthermore, if the determinant of a matrix is negative it is definitely not positive semidefinite (since positive semidefinite eigenvalues are non-negative and the determinant is the product of the eigenvalues), but we cannot say anything if the determinant is positive.

### 5.1. Common Kernels and Rules

- Linear kernel:  $\kappa(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y} + c, c \geq 0$
- Polynomial Kernel  $\kappa(\mathbf{x}, \mathbf{y}) = (a\mathbf{x}^\top \mathbf{y} + c)^d, a, c, d \geq 0$
- Gaussian Kernel  $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}), \sigma > 0$
- Exponential Kernel  $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x} - \mathbf{y}\|_2}{\sigma}), \sigma > 0$
- If  $\kappa_1, \kappa_2$  are Kernels and  $c \geq 0$ , then  $c\kappa_1, c + \kappa_1, \kappa_1 + \kappa_2$  and  $\kappa_1 \cdot \kappa_2$  are kernels as well

Notes:

## 6. Principal Component Analysis

Assuming distribution of unlabeled raw data is concentrated around some lower-dimensional plane and a projection onto this plane captures most of the variance of the dataset. PCA uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.

### 6.1. Geometric Interpretation

- Imagining that  $x$  is a vector in  $p$  dimensions, and  $U_k$  describes a lower dimensional subspace, then the coordinates of  $x$  projected onto  $U_k$  is found by first obtaining the scores of  $x$  along  $U_k$  and then taking the linear combination of these scores and the  $U_k$  vectors:  $\pi_U(x) = U_k U_k^T x$ ,  $U_k \in \mathbb{R}^{p \times k}$
- We want to find a  $U_k$  that captures most of the variance in the dataset, which is equivalent to finding a hyperplane where the difference between the original data points and their projections are minimized, thus forming our problem for optimization.

$$J(U_k) = \sum_{i=1}^n ||x_i - U_k U_k^T x_i||_2^2$$

- The dataset can be represented as the dot product of an orthonormal matrix, a diagonal matrix and another orthonormal matrix via singular value decomposition:  $X = UDV^T$ .

The first  $k$  columns of the first orthonormal matrix via SVD  $U_k$  minimizes the difference between data points and their projections, and the corresponding covariance matrix of the scores is diagonal (= features are uncorrelated).

- The *empirical covariance matrix* or *Score Matrix* of the reduced  $k$  variables  $S = U_k^T X$  is diagonal, i.e. the features extracted by the PCA are uncorrelated.  $[cov(S) = \frac{1}{n} S S^T = \frac{1}{n} U^T X X^T U]$   
 $s_{ij}$ :  $j$ -th score for the  $i$ -th principal component
- $U_k$  is the Loading Matrix give by  $k$  left-singular vectors of the SVD of  $\tilde{X}$ .
- $S$  can directly be obtained by the SVD with  $S = D_k V_k^T$ . As  $D_k$  is diagonal and  $V_k^T$  only has  $k$ -rows, this only requires  $nk$  operations instead of non-feasible  $(2p - 1)nk$  operations.
- Dimensions:**  
 $X \in \mathbb{R}^{m \times n}$ ,  $U \in \mathbb{R}^{m \times m}$ ,  $D \in \mathbb{R}^{m \times n}$ ,  $V^T \in \mathbb{R}^{n \times n}$ ,  
 $U_k \in \mathbb{R}^{k \times n}$ ,  $X_k \in \mathbb{R}^{k \times n}$ ,  $m$ : Dimension variable,  $n$ : samples
- The eigenvectors of  $XX^T$  are the loadings of  $X$  only if  $X$  is centered.

### 6.2. Proof

- Theorem 4.1 is proved by first reframing the minimization problem (4.2) into a maximization problem, reframing this again into a problem involving the trace of a rank  $k$  projection matrix, then noting that the maximum value along the diagonals is 1 (4.5) of which subbing in  $U_k$  achieves due to its resulting in a identity matrix. We can show that the scores are uncorrelated by showing that the dot product between the scores and its transpose scaled by  $1/n$  is a diagonal matrix using the SVD. (4.6)
- If we have the product between two matrices where the first matrix is comprised of a diagonal matrix as well as a zero matrix, then we can get rid of the zero columns in the first matrix and also get rid of the corresponding rows in the left matrix.
- A computationally inexpensive way of obtaining the scores using the singular values and right singular vector can be derived by substituting the SVD into the formula for the calculating the scores. (4.7)

### 6.3. Statistical Interpretation

The idea behind the statistical interpretation is that covariance matrices are symmetric positive semidefinite, and therefore there exists some matrix  $U$  which can diagonalize it (clearly illustrated using the eigenvalue decomposition). i.e.  $Y = U^T X$  with  $U^T U = I_p$ , such that the covariance matrix of  $Y$  is diagonal and the variances decrease. In the same way:

$$\begin{aligned} D &= U^T var(X) U = \mathbb{E}[U^T (x - \mu_x)(x - \mu_x)^T U] \\ &= \mathbb{E}[(U^T X - U^T \mu_x)(U^T X - U^T \mu_x)^T] \\ &= \mathbb{E}[(Y - \mu_y)(Y - \mu_y)^T] = var(Y) \end{aligned}$$

### 6.4. Error Model Interpretation

The projection is chosen in a way that the error measured with the Squared Euclidean Norm is minimized. The minimized noise is *White Gaussian*.

- The error model interpretation is that the samples  $X$  is obtained via the addition between a signal matrix  $L$  (*clean data*) and a noise matrix  $N$ , where  $L$  is lying on a lower dimensional subspace than  $X$ , and the goal is to try and find  $L$ .  $[X = L + N]$
- Setting  $L$  equal to the SVD ( $p$  dimensions with  $k$ ;  $rank(L) \leq k$ ) minimizes the Frobenius norm of the difference between  $X$  and  $L$  given that the dimension  $k$  of the subspace which  $L$  lies in is known beforehand.  
 $\hat{L} = \arg \min_{rank(L)} ||X - L||_F = U_k \text{diag}(d_1, \dots, d_k) V_k^T$

### 6.5. Relation to Autoencoders

- Autoencoders map an input to low dimensional space using a function  $f$ , which is then mapped back to a higher dimensional space using a function  $g$  and can approximate the input  $g \circ f(x_i) \approx x_i$
- Letting  $f$  and  $g$  be linear mappings represented by matrices, the reconstruction error  $J(W, V) = \sum_{i=1}^n ||x_i - wVx_i||_2^2$  is minimised by setting  $f$  as the transpose of  $U_k$  and  $g$  as  $U_k$ .

Let  $U_k$  be the first  $k$  left singular vectors of observation matrix  $X$ , the  $V = U_k^T$  and  $W = U_k$  minimize the reconstruction error of the linear autoencoder.

- This is proven by noting that in Theorem 4.2 which has a problem of the same form, the error is minimized using the projections of the points on the subspace of interest, thus leading to  $f$  and  $g$  as the transpose of  $U_k$  and  $U_k$  itself respectively as they cause this projection.

#### Notes:

- Low dimensional subspace  $\hat{=}$  low rank

## 7. Kernel PCA

### 7.1. Linear PCA with Inner Products

- Scores of the input data can be calculated using singular values and the right eigenvector (4.7), which we can obtain using the eigenvalue decomposition of the Gram-Matrix (7.2) and therefore giving us a way to calculate linear PCA scores using the inner product. If we forgot to center the data, use the centred Gram-Matrix instead (7.8).  $[V$  is orthogonal;  $\Sigma^T \Sigma$  is diagonal]
- If the input data has not been centred beforehand, we can find the Gram-Matrix for the centred data by noting how input data is usually centred (7.6), factoring out the input matrix (7.7) and then using inner product to compute the Gram-Matrix corresponding to centred input. (7.8)
- Scores for new observations are usually calculated by projecting them onto  $U_k$ , but if we want to reformulate this to be in terms of the singular value and right singular vectors we can isolate  $U$  in the SVD and retain only  $k$  columns (7.4a) to arrive at the formulation. (7.4)
- If we want to calculate scores for new observations considering centering, we take the formula for raw new observations (7.4), replace the input matrix with a centred input matrix and the observation with an approximated centred observation, resulting in a nice formula involving the Gram-Matrix. (7.9)

EVD of Gram Matrix (7.2):  $K = V \Sigma^T \Sigma V^T$   
Centered Gram Matrix (7.8):  $\tilde{K} = H K H$  with  
 $H = (I_n - \frac{1}{n} \mathbb{1}_n \mathbb{1}_n^T)$

Computing scores for new sample  $y$ :  $U_k^T(\tilde{y}) = \Sigma_k^{-1} V_k^T \tilde{k}_y$ ,  
 $\tilde{k}_y = H k_y - \frac{1}{n} H K \mathbb{1}_n$

- With  $V_k$  and  $\Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k)$  are the first  $k$  eigenvectors and eigenvalues of  $K$  we can compute the principal components  $U_k^T y$  by only using inner products.

$U_k^T = \Sigma_k^{-1} V_k^T X^T$  and for a new sample  $y$ :  
 $U_k^T y = \Sigma_k^{-1} V_k^T X^T y = \Sigma_k^{-1} V_k^T [x_1^T y \dots x_n^T y]^T$   
This only requires the inner product  $x_n^T y = k_y$

### 7.2. Transition to Kernel PCA

- Kernel PCA requires us to first compute the Gram-Matrix via the kernel (6.1). This can then be substituted for the Gram-Matrix via inner product at various places, such as during centering (7.8) and when performing eigenvalue decomposition (7.2) to calculate matrices for scores (4.7).
- The scores for new observations via Kernel PCA is based off the Linear PCA version (7.4), and can be obtained by replacing the linear observation inner products with the kernel observation inner products (as well as ensuring that singular values and right singular vectors obtained from eigenvector decomposition of Gram-Matrix via kernel (7.2)).
- Similarly, when considering centering, note the Linear PCA version of scores which consider centering (7.9) and replace the linear Gram-Matrix with the kernel Gram-Matrix as well as the linear observation inner products with the kernel observation inner products. (7.12)

Instead of replacing the inner product  $x^T y$  by  $\langle \phi(x), \phi(y) \rangle$ , we substitute  $x^T y \rightarrow \kappa(x, y)$ :

$$\begin{aligned} k^{new} &= [\kappa(x_1, y) \dots \kappa(x_n, y)]^T \\ k_{cent}^{new} &= H k^{new} - \frac{1}{n} H K \mathbb{1}_n \end{aligned}$$

### Kernel Principal Component Analysis

- for training set  $X = [x_1 \dots x_n]$ ,  $x_i \in \mathbb{R}^p$
- Find suitable Kernel function  $\kappa(\cdot)$  and compute Gram Matrix

$$K = \begin{bmatrix} \kappa(x_1, x_1) & \dots & \kappa(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \kappa(x_n, x_1) & \dots & \kappa(x_n, x_n) \end{bmatrix}$$

- Compute centered Gram Matrix:  $\tilde{K} = H K H$  with  $H = I_n - \frac{1}{n} \mathbb{1}_n \mathbb{1}_n^T$
- Compute Eigenvalue Decomposition:  $\tilde{K} = V \Lambda V^T$ . Because  $K$  is positive semi-definite and therefore the diagonal entries of  $\Lambda$  are non-negative, we write  $\Lambda = \Sigma^2 = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$
- Reduced Matrices:  $\Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k)$ ,  $\in \mathbb{R}^{k \times k}$  and  $V_k \in \mathbb{R}^{n \times k}$
- Reduced Training Data:  $S = \Sigma_k V_k^T$
- For new datapoint  $y \in \mathbb{R}^p$ , compute new components:

$$\begin{aligned} s_{new} &= \Sigma_k^{-1} V_k^T k_{cent}^{new} \\ k_{cent}^{new} &= H k^{new} - \frac{1}{n} H K \mathbb{1}_n \\ k^{new} &= [\kappa(x_1, y) \dots \kappa(x_n, y)]^T \end{aligned}$$

#### Notes:

## 8. Feedforward Neural Networks

### 8.1. Definitions and Motivation

The power behind an FNN is its ability to minimize any kind of expected prediction error, that is, find model parameters that minimize expected loss (5.1). Feed forward neural networks are concatenations of linear functions (5.2) with activation functions (e.g. 5.3), which can be denoted in vector form. This results in the generalized form of an FNN. (5.5)  
With  $p/N$  being Input/Output dimensions and  $l$  layers.  $N$  depends on loss function (Output of FNN is input of Loss). A Deep FNN has  $l > 3$ .

Given function class  $\mathcal{F}$  described by set of  $N$  parameters  $\Theta \in \mathbb{R}^N$ :  
Solve:  $\hat{\Theta} = \arg \min_{\Theta \in \mathbb{R}^N} \frac{1}{n} \sum_i L(f_{\Theta}(\mathbf{x}_i))$  (5.1)  
 $\varphi_{\mathbf{W}}: \mathbb{R}^p \rightarrow \mathbb{R}^m, \mathbf{h} \rightarrow \mathbf{W}\mathbf{h}$  (5.2)  
 $\sigma(t) = \max\{0, t\}$  (ReLU, 5.3)  
 $f: \mathbb{R}^p \rightarrow \mathbb{R}^o: \sigma_l \circ \varphi_{\mathbf{W}_l} \circ \dots \circ \sigma_1 \circ \varphi_{\mathbf{W}_1}(\mathbf{x})$  (5.5)

### 8.2. Training FNN

The gradients can be calculated via chain rule through each subsequent layer, or in this case due to vectors via the Jacobian. From (5.5) the key stages include matrix multiplication (5.2) and nonlinear activation (5.3). The corresponding Jacobians w.r.t. inputs in each stage is seen in (5.6) and (5.7) respectively. However, for the Jacobian of matrix multiplication w.r.t weights, the weights need to be reshaped into a vector (5.9) which will then allow us to formulate the Jacobian (5.10). Then the gradients w.r.t loss (5.11) is simply (5.31a). The gradient can then be used to update the weights for the next time step, making sure to apply an inverse reshaping of the gradient vector into a matrix (5.14).

$\mathbf{J}_{\mathbf{W}_g}(\mathbf{x}) = \mathbf{W} \cdot \mathbf{J}_g(\mathbf{x})$  (5.6)  
 $\mathbf{J}_{\sigma}(\mathbf{x}) = \text{diag}(\sigma'(x_1), \dots, \sigma'(x_m))$  (5.7)  
 $\mathbf{J}_{mult}(\mathbf{x}) = \text{diag}(\mathbf{x}_1^T, \dots, \mathbf{x}_n^T) \in \mathbb{R}^{m \times mn}$  (5.10), doesn't depend on  $\mathbf{W}$   
With  $\mathbf{h}_j = \sigma_j \circ \varphi_{\mathbf{W}_j} \circ \dots \circ \sigma_1 \circ \varphi_{\mathbf{W}_1}(\mathbf{x})$  the  $j$ -th output  
We get for the output-layer (L):  
 $\frac{\partial F}{\partial \mathbf{W}_L} = \mathbf{J}_L(\mathbf{h}_L) \cdot \mathbf{J}_{\sigma_L}(\mathbf{W}_L \mathbf{h}_{L-1}) \cdot \mathbf{J}_{mult}(\mathbf{h}_{L-1})$   
 $\frac{\partial F}{\partial \mathbf{W}_j} = \mathbf{J}_L(\mathbf{h}_L) \cdot \mathbf{J}_{\sigma_L}(\mathbf{W}_L \mathbf{h}_{L-1}) \cdot \dots \cdot \mathbf{J}_{\sigma_j}(\mathbf{W}_j \mathbf{h}_{j-1}) \cdot \mathbf{J}_{mult}(\mathbf{h}_{j-1})$  for the  $j$ -th layer

### 8.3. Multiclass Classification with FNNs

Considers the problem of assigning an input  $\mathcal{X} \in \mathbb{R}^p$  to one out of multiple, say  $C$  classes.

- Model problem with random variable  $(\mathcal{X}, \mathcal{Y})$ , with  $\mathcal{X} \in \mathbb{R}^p$  and  $\mathcal{Y} \in \{\mathbf{e}_1, \dots, \mathbf{e}_C\}$  as one of the standard basis vectors in  $\mathbb{R}^C$ .
- Realization  $\mathbf{y} = \mathbf{e}_c$  means that the event *belongs to class c* is true.  $\rightarrow$  This modelling of the output variable is also known as *one-hot-encoding*.
- Idea:  $\mathcal{X}$  is the input to the FNN and output vector  $\mathbf{h}_l$  in  $\mathbb{R}^C$  approximates the probability of the class distribution given  $\mathcal{X}$ .

$$\mathbf{h}_l \approx \begin{bmatrix} Pr(\mathcal{Y} = \mathbf{e}_1 | \mathcal{X} = \mathbf{x}) \\ \vdots \\ Pr(\mathcal{Y} = \mathbf{e}_C | \mathcal{X} = \mathbf{x}) \end{bmatrix}$$

- Thus last activation function  $\sigma_l$  outputs probability distribution for the  $C$  classes. All entries add up to 1.
- Example for output function is the **softmax**, given  $\forall c \in 1, \dots, C$  by

$$\sigma: \mathbb{R}^C \rightarrow \mathbb{R}^C, \begin{bmatrix} a_1 \\ \vdots \\ a_C \end{bmatrix} \rightarrow \begin{bmatrix} S_1 \\ \vdots \\ S_C \end{bmatrix}, S_c = \frac{\exp a_c}{\sum_c \exp a_c}$$

- Derivative/Jacobi matrix  $J$  of softmax for arbitrary  $i$  and  $j$ :

$$J_{i,j} = \begin{cases} S_i(1 - S_j) & , i = j \\ -S_j S_i & , i \neq j \end{cases}$$

## 9. Support Vector Machines

### 9.1. Geometry

Idea: Find hyperplane to divide  $x, y$  into two subspaces

- For some vector  $w$ , an affine hyperplane is described as the set of points  $x$  that can be projected onto  $w$  and then shifted by some constant  $b$  to equal zero. (8.1)
- The vector  $w$  is normal to the hyperplane (8.1) since a vector of any direction in the hyperplane projected onto  $w$  has a magnitude of zero. (8.2)
- The signed distance from any point to a hyperplane is defined as in (8.3), and can be interpreted as the distance from the hyperplane to the point  $x$ . It is formulated by first projecting  $x$  onto  $w$  and then centering the number line at  $b$  to obtain the signed value of  $x$  in the decision perspective (denote this as the decision distance). Then this value is divided by the magnitude of  $w$  to obtain the signed distance. Note that  $-b$  denotes the decision distance to the origin, and the magnitude of  $w$  is a scaling factor used to shift between the decision distance and signed distance.

$$\begin{aligned} \text{Hyperplane: } \mathcal{H}_{w,b} &= \{x \in \mathbb{R}^n | w^T x - b = 0\} \quad (8.1) \\ w &\text{ is normal to } \mathcal{H}: w^T (x_1 - x_2) = 0 \quad (8.2) \\ \text{Signed Distance: } \delta(x, \mathcal{H}_{w,b}) &= \frac{w^T x - b}{\|w\|} \quad (8.3) \\ \text{Margin: } \mathcal{H}_+ &= \{x \in \mathbb{R}^n | w^T x - b = 1\} \\ \mathcal{H}_- &= \{x \in \mathbb{R}^n | w^T x - b = -1\} \quad (8.8-9) \end{aligned}$$

- The decision perspective can be imagined as the perspective where classification decisions are made, all depending on the decision distance.
- The positive margin is defined as the hyperplane with a decision distance of 1 (8.8), and the negative margin is defined as the hyperplane with a decision distance of -1 (8.9).
- Decision boundary of SVM is always a hyperplane (i.e. an affine space) in  $\mathbb{R}^p$ .
- Retraining with a new training sample affects the decision boundary.

### 9.2. Basic Linear SVM

- The aim of Linear SVM is to linearly separate data using the margins of the hyperplanes, such that all positively labelled data is bounded by the positive margin and all negatively labelled data is bounded by the negative margin. In other words, all positively labelled data has a decision distance greater than or equal to 1, and all negatively labelled data has a decision distance less than or equal to -1, resulting in a separation constraint for  $w$  and  $b$ . (8.12)
- We define the best  $w$  and  $b$  as the ones that have the widest margin, i.e. the widest distance between the hyperplanes. This width is calculated by the sum of the signed distance from the positive margin to the hyperplane and the signed distance from the hyperplane to the negative margin (8.13), which we can then maximize subject to our separation constraint (8.12).
- Flipping around the numerator and denominator of the width, this becomes a minimization problem (8.14).

$$\begin{aligned} \text{Conditions: } w^T x_i - b &\geq +1 \text{ for } y_i = 1 \text{ and } \\ &w^T x_i - b \leq -1 \text{ for } y_i = -1 \\ \text{Compact: } y_i \cdot (w^T x_i - b) &\geq 1 \text{ for all } i = 1, \dots, N \quad (8.12) \\ \text{Minimization (8.14):} \\ \min_w \frac{1}{2} \|w\|_2^2 \text{ s.t. } &y_i (w^T x_i - b) \geq 1 \text{ for all } i = 1, \dots, N \end{aligned}$$

### 9.3. Karush-Kuhn Tucker Conditions

- The KKT conditions say that a minimisation problem with a set of equality and inequality constraints (8.15) can be reformulated as a Lagrangian primal (8.16) with some KKT conditions (Theorem 8.2), such that for convex objective functions with a convex feasible region (like that defined 8.14) minimising the primal (8.16) is equivalent to minimising the original problem (8.15).

For an optimization problem with  $\mathcal{E}$  and  $\mathcal{I}$  as equality and inequality constraints (8.15):  
 $\min_{z \in \mathbb{R}^n} f(z) \text{ s.t. } c_i(z) = 0 \forall i \in \mathcal{E} \text{ and s.t. } c_j(z) \geq 0 \forall j \in \mathcal{I}$   
Lagrange Function (8.16):  $L(z, \lambda) = f(z) - \sum_{i \in \mathcal{I} \cup \mathcal{E}} \lambda_i c_i(z)$   
**Karush-Kuhn-Tucker Conditions**  
For  $z^*$  as a solution to 8.15, there exists a Lagrange multiplier  $\lambda^*$  such that (8.17-8.21):  
 $-\nabla_z L(z^*, \lambda^*) = 0$   
 $c_i(z^*) = 0 \quad \forall i \in \mathcal{E}$   
 $c_i(z^*) \geq 0 \quad \forall i \in \mathcal{I}$   
 $\lambda_i^* \geq 0 \quad \forall i \in \mathcal{I}$   
 $-\lambda_i^* c_i(z^*) = 0 \quad \forall i \in \mathcal{I} \cup \mathcal{E}$

- The convex primal function (8.16) can then be reformulated as a concave dual function (8.22) by taking the infimum of the primal function. The infimum is the set of points along which the function is minimized w.r.t. the non-Lagrangian multiplier variables. Then by maximizing the primal subject to its constraints (8.22a) we obtain a lower bound for the solution to the primal (weak duality) or when some conditions are satisfied (such as in SVM) this coincides with the solution to the primal (strong duality).

### 9.4. Linear SVM via Lagrangian Duality

- SVM via Lagrangian duality follows the process specified above. The original problem is as in (8.14), the Lagrangian primal is as in (8.23 to 8.25) and the KKT conditions are as in (8.17 to 8.21).

$$\begin{aligned} \text{Problem: } \min_{w,b,\lambda \geq 0} L(w,b,\lambda) \\ L(w,b,\lambda) &= \frac{1}{2} \|w\|^2 - \sum_i \lambda_i y_i (w^T x_i - b) + \lambda_i \quad (8.25) \\ \nabla_{(w,b)} L(w,b,\lambda) &= \begin{bmatrix} w - \sum_i \lambda_i y_i x_i \\ -\sum_i \lambda_i y_i \end{bmatrix} \\ \text{KKT Conditions:} \\ w^* - \sum_i \lambda_i^* y_i x_i &= 0 \quad (8.27), \quad \sum_i \lambda_i^* y_i = 0 \quad (8.28) \\ \lambda_i^* (y_i (w^{*T} x_i - b^*) - 1) &= 0 \quad (8.29) \\ \text{Returns (8.30):} \\ \min_{w,b,\lambda \geq 0} L(w,b,\lambda) &= L_D(\lambda) = \\ \sum_i \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j &\text{ s.t. } \lambda_i \geq 0, \sum_i \lambda_i y_i = 0 \end{aligned}$$

- Then the dual function needs to be calculated by taking the infimum (*greatest element*) of (8.25), which is accomplished by substituting (8.27) and then (8.28) into (8.25), resulting in the dual form (8.30).
- Maximising this function w.r.t. its constraints is then the dual problem for the SVM. (8.30)
- Problem is strictly convex, therefore solution is unique
- If  $\lambda_i \neq 0$ , then  $x_i$  is support vector (Lies in or on margin)
- Only works, if classes can be linearly separated  
Else: Kernel/Soft Margin SVM required
- SVM works better when (nearly) separable than Linear Regression and is preferred when using kernels.

### 9.5. Soft Margin Linear SVM

Allows for some wrongly assigned data samples and searches for a trade-off between a large margin and a degree of misclassification.

- Misclassification is quantified by set of  $N$  additional variables  $\xi_i$  that are assumed to be non-negative.  
 $y_i (w^T x_i - b) \geq 1 - \xi_i$ , for all  $i = 1, \dots, N$
- Optimization problem changes to  
 $\min_{w,b,\xi} \frac{1}{2} \|w\|_2^2 + c \sum_{i=1}^N \xi_i$   
s.t.  $y_i (w^T x_i - b) \geq 1 - \xi_i$  and  $\xi_i \geq 0 \quad \forall i$   
[  $c > 0$  and weighs between large margin and misclassification  
 $\rightarrow$  larger  $c$ , the violation of the separation rule is punished more.]

### 9.6. Kernel SVM

- Dual Form:  
 $\max_{\lambda} \sum_i \lambda_i - \frac{1}{2} \lambda^T H \lambda \text{ s.t. } 0 \leq \lambda_i \leq c, \sum_i \lambda_i y_i = 0$   
with  $h_{ij} = y_i y_j x_i^T x_j$  which changes to  $h_{ij} = y_i y_j \kappa(x_i, x_j)$  using a Kernel function.
- Using the KKT conditions, we can compute  
 $b = \frac{1}{N_{Supp}} \sum_{i \in Supp} (\sum_{j \in Supp} \lambda_j y_j \kappa(x_i, x_j)) - y_i$
- $Supp$  denotes all  $i$  for which  $0 < \lambda_i \leq c$ .

## 10. Useful Facts

- The matrix resulting from the dot product of between an  $\mathbb{R}^n \times p$  matrix and  $\mathbb{R}^{p \times n}$  has at most a rank of  $p$ .
- Examples for Loss-functions:  
**Autoencoder:**  $L(f(x)) = \|f(x) - X\|_2^2$   
reconstruct samples of  $x$   
 $f$ : concatenation of encoders and decoder function  
**Regression:**  $L(f(x), Y) = \|f(x) - Y\|_2^2$   
best  $f \in F$  that minimizes expectation of  $L$
- Unsupervised Learning:** data does not have to be labeled (by supervisor) before data is employed for learning algorithm (e.g. PCA).
- Supervised Learning:** data is labeled and thus the model is trained  $\rightarrow$  Classification vs. Regression
- Regression:** predicting trends using previous labeled data
- $cov(x_1, x_2) = \frac{1}{n} X X^T = \frac{(x_{11}x_{21} + x_{12}x_{22} + \dots)}{n \cdot \|x_1\| \cdot \|x_2\|}$
- Data Scales:**  
*Qualitative:* Nominal ( $=, \neq \rightarrow$  Gender or Names), Ordinal ( $\geq, \leq, >, < \rightarrow$  Happiness)  
*Quantitative:* Interval ( $x, - \rightarrow$  2019 or 18°C), Ratio ( $*, \div \rightarrow$  Distance or Area)
- A function  $f$  is called a **Distance Measure** if  $\forall x, y \in \mathbb{R}^p$ :  
*Symmetry:*  $d(x, y) = d(y, x)$   
*Positive Definiteness:*  $d(x, y) = 0 \Leftrightarrow x = y; d(x, y) \geq 0$   
*Triangular Irregularity:*  $d(x, z) \leq d(x, y) + d(y, z)$
- Taxicab/Manhattan Norm:**  $\|x\|_1 = \sum_{i=1}^n x_i$   
**Euclidean Norm:**  $\|x\|_2 = \sqrt{x_1^2 + \dots + x_n^2}$   
**Maximum Norm:**  $\|x\|_\infty = \max(|x_1|, \dots, |x_n|)$   
**Inner Product Norm:**  $\|x\|_A = \sqrt{x A x^T}$  [Frobenius:  $A = 1$ ]  
*Reminder:*  $\|q \cdot x\| = |q| \cdot \|x\|$  and  $\|x + y\| \leq \|x\| + \|y\|$
- Lagrange Method:**  
1.  $L(x, y, \lambda) = f(x, y) + \lambda \cdot g(x, y)$   
( $g(x, y)$  is the conditional equation set to zero)  
2.  $\nabla L \stackrel{!}{=} 0$  (calculate  $L_x(\text{I})$ ,  $L_y(\text{II})$  and  $L_\lambda \stackrel{!}{=} 0$  (III))  
3. Eliminate  $\lambda$  from  $L_x$  &  $L_y$ , (IV):  $\frac{df(x,y)}{dx} \cdot y - \frac{df(x,y)}{dy} \cdot x \stackrel{!}{=} 0$   
4. Solve system of equations (III) and (IV) to find all  $x$  and  $y$ .
- Training and test sets normally get divided into:  
*Machine Learning:* 60/20/20  
*Deep Learning:* 98/1/1

## 11. Homework and Assignments

$p_X(X_1, X_2)$	$X_2 = 0$	$X_2 = 1$
$X_1 = 0$	$p_X(0, 0)$	$p_X(0, 1)$
$X_2 = 1$	$p_X(1, 0)$	$p_X(1, 1)$

Calculate Covariance Matrix from Table	
1. Calculate Means for $X_1$ and $X_2$ : $\mu_1, \mu_2$	
2. Create X-Matrix, e.g.:	$\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$
3. Create p-Matrix:	$\begin{bmatrix} p_X(0, 0) & 0 & 0 & 0 \\ 0 & p_X(1, 0) & 0 & 0 \\ 0 & 0 & p_X(0, 1) & 0 \\ 0 & 0 & 0 & p_X(1, 1) \end{bmatrix}$
4. Calculate Covariance: $Cov = XpX^T$	

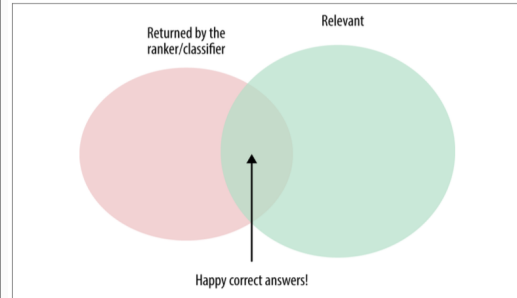
### 11.1. Classification Analysis

**False Positive:** Indicates presence of a condition during a test, but condition is actually **not** present.

**False Negative:** Indicates **no** presence of a condition, actually **is** present.

**True Positive:** Indicates presence of a condition, actually **is** present.

**True Negative:** Indicates **no** presence of a condition, actually **not** present.



• **Accuracy:**  $\frac{\text{correct predictions}}{\text{total data points}}$

Drawbacks: Cost of misclassification might differ for different tasks & we have a lot of data for one class but not the other (*Skewed Data*).

• **Precision:**  $\frac{\text{happy correct answers}}{\text{total items returned by ranker}} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$

• **Recall:**  $\frac{\text{happy correct answers}}{\text{total relevant items}} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$

• **F1-Score:**  $F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$

• **NDCG (Normalized Discounted Cumulative Gain):**

Similar to a search engine, top few answers in a query matter more. **Related metrics:** Cumulative gain sums up the relevance of the top  $k$  items. Discounted cumulative gain (DCG) discounts items that are further down the list. Normalized discounted cumulative gain is a normalized version of DCG which divides the DCG by the perfect ideal DCG score with ideal ordering, so that the normalized score always lies between 0.0 and 1.0. For a position  $p$ :

$$DCG_p = \sum_{i=1}^p \frac{2^{rel_i - 1}}{\log_2(i+1)} \quad [\text{for IDCG } p = |REL|]$$

#### 11.1.1. ROC Curve

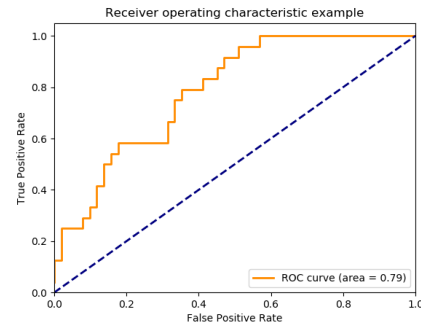
The ROC curve is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied. It shows:

- The relationship between sensitivity and specificity. For example, a decrease in sensitivity results in an increase in specificity.
- Test accuracy; the closer the graph is to the top and left-hand borders, the more accurate the test. Likewise, the closer the graph to the diagonal, the less accurate the test. A perfect test would go straight from zero up the the top-left corner and then straight across the horizontal.

- The likelihood ratio; given by the derivative at any particular cutpoint.

• **False Positive Rate:**  $\frac{\text{False Positives}}{\text{False Positives} + \text{True Negatives}}$

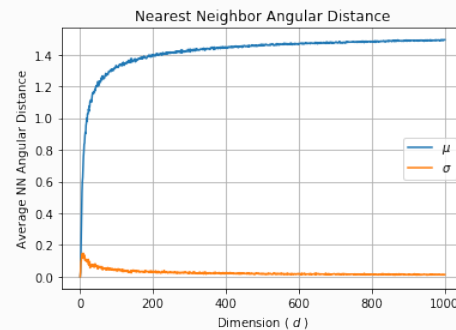
**True Positive Rate:**  $\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$



The **Area under the Curve (AUC)** is the integral of the ROC Curve and gives a measure how good the classifier is. An area under the ROC curve of 0.8, for example, means that a randomly selected case from the group with the target equals 1 has a score larger than that for a randomly chosen case from the group with the target equals 0 in 80% of the time. When a classifier cannot distinguish between the two groups, the area will be equal to 0.5 (the ROC curve will coincide with the diagonal). When there is a perfect separation of the two groups, i.e., no overlapping of the distributions, the area under the ROC curve reaches to 1 (the ROC curve will reach the upper left corner of the plot).

### 11.2. Curse of Dimensionality

- The angular distance between 2 randomly sampled vectors increases with dimension  $d$  of the sample space.
- Convergence to  $\frac{\pi}{2}$  implies that two randomly sampled vectors are orthogonal to each other in  $d$ -dimensional apsce for  $d \gg n$ .
- Convergence to  $\frac{\pi}{2}$  also implies that most samples are concentrated in the 'corners' of the  $d$ -dimensional cube  $[-1, 1]^d$ , i.e. in high dimension, the corners occupy most of the space.
- This convergence also means that 2 randomly sampled vectors are increasingly equidistant (in terms of angular distance) from their respective nearest neighbors in high dimensional space.
- Because the samples are increasingly equidistant from each other, this means that distance-based classifiers (e.g. k-Nearest Neighbors) cannot be used on such data in high-dimensional space.
- Increasing the sample size  $n$  decreases the average angular distance between neighbouring vectors in a  $d$ -dimensional feature space. The rate of decrease, however, decreases with increasing  $n$ .



### 11.3. Logistic Regression

- With big datasets, standard gradient descent could lead to Memory Error
- Use *stochastic gradient descent* instead: Train over epochs instead:
- Each epoch, the training set is divided randomly into equal size subsets (=minibatch). Then the gradient of each subset is calculated and applied only to the samples in the subset
- A epoch is finished when the gradient step was performed on each subset

### 11.4. Principal Component Analysis

Removing the first  $n$  columns from  $U_k$  can have different effects on classification:

- **Decreased Error Rate:** This may be because even though the first  $n$  components capture more variance in the samples, perhaps the other components are better at separating samples by labels, allowing KNN to correctly classify samples (Subset 1 in top plot and second plot)
- **No Effect on Error Rate:** This may be because the first three principal components are as good at separating samples by labels compared to other principal components (Subsets 2+3 in top plot and third plot)
- **Increase Error Rate:** This may be because the first three principal components are better at separating samples by labels compared to the other principal components (Subset 4 top plot and bottom plot)

#### 11.4.1. How to choose $k$ ?

Assuming that  $\mathbf{X} \in \mathbb{R}^{p \times N}$  is the centered data matrix and  $\mathbf{P} = \mathbf{U}_k \mathbf{U}_k^T$  is the projector onto the  $k$ -dimensional principal subspace, the dimension  $k$  is chosen such that the fraction of overall energy contained in the projection error does not exceed  $\epsilon$ , i.e.

$$\frac{\|\mathbf{X} - \mathbf{P}\mathbf{X}\|_F^2}{\|\mathbf{X}\|_F^2} = \frac{\sum_{i=1}^M \|\mathbf{x}_i - \mathbf{P}\mathbf{x}_i\|^2}{\sum_{i=1}^N \|\mathbf{x}_i\|^2} \leq \epsilon,$$

where  $\epsilon$  is usually chosen to be between 0.01 and 0.2. Energy is not always the best way to measure useful information, e.g. when images differ in brightness (=No use full information)

### 11.5. Exam Winter Term 18/19

The questions were completely different to what was discussed in the lecture, the assignments and the tutorial. They focused on logical thinking and simple calculations. Questions had a structure like such:

"These six different  $x_i$  and  $y_i$  are given. Which of the following  $\mathbf{w}$  and  $\lambda$  are correct?"

→ Solution: Use the KKT Conditions with each possible  $\mathbf{w}^*$  and  $\lambda^*$  to see which one has a plausible solution.

"Which of the following could (not) be a Kernel function?"

→ Solution: Check if  $K$  is positive semidefinite, if  $\kappa$  is symmetric etc.

Most questions could be solved by using the method of elimination and almost all the time the information necessary to solve the question was in the description of the question (e.g. they stated the KKT Conditions).