



Information Retrieval in High Dimensional Data

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 $\mu = E[X] \in \mathbb{R}^p$
 Estimate mean: $\hat{\mu}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}$ (Sum row wise!)
 Variance $\Sigma = \text{Var}(X) = E[(X - \mu)(X - \mu)^T] \in \mathbb{R}^{p \times p}$
 Σ is symmetric and positive semidefinite: $x^T \Sigma x \geq 0 \forall x$ or $\Sigma \geq 0$
 Centered Observation Matrix: $\hat{X}_{ij} = x_{ij} - \hat{\mu}_i$ Transpose Matrix: $(XY)^T = Y^T X^T$

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: $p_{Y|X=x}(y) = \frac{p_{X,Y}(x,y)}{p_X(x)}$ Expectation Value: $E[X] = \int_{\mathbb{R}^p} x p_X(x) dx = \mu_X$
 (Co)Variance: $\text{Var}[X] = E[(X - \mu_X)(X - \mu_X)^T]$

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, conditional median is best

2.2. Decision Making

- Global Methods:** Find best explicit global function $\hat{f} = \text{argmin}_{f \in \mathcal{F}} EPE(f)$, incorporate all points reduce complexity by learning parametrize decision function
- Local Methods:** Find best local value \hat{c} of given realization x
 $\hat{c} = \text{argmin}_{c \in \mathbb{R}} E_{Y|X=x} L(Y, c)$, only samples in region of interest

2.3. Curse of Dimensionality

With increasing dimension p :

- Noise increases (accumulates over dimensions)
- Number of observations for same accuracy increases exponentially
- Empty Space Phenomenon:** High dimensional spaces are sparse
- Bigger space = more samples in the tail of the distribution
- Samples are equidistant to each other

Therefore, difficult to estimate underlying distribution

2.4. 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. Logistic Regression

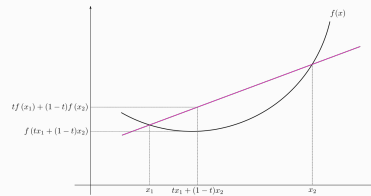
3.1. Binary Classification

Input Data: $X \in \mathbb{R}^p$, Output Variables: $Y \in \{-1, +1\}$
 $L_{0,1}(Y, f(X)) = 1$ if $Y \cdot \text{sign}(f(X)) \leq 0$, numerically infeasible

3.2. Convexity

in convex set $C \subset \mathbb{R}$, f is convex, if $t f(x_2) + (1-t)f(x_1) \leq f(t x_2 + (1-t)x_1) \forall x_1, x_2 \in C, t \in [0, 1]$

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



3.3. Logistic Regression

Choose $f(x) = w^T x + b$ and log-loss $l(t) = \log(1 + e^{-t})$:

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

$P(Y = y|x) = \frac{1}{1 + \exp(-y(w^T x + b))}$, Prob. that y is correct

Find w^*, b^* by gradient descent
 Classify new sample: $\text{sign}(w^{*T} x_{new} + b^*)$

3.4. Overfitting

For linearly separable non empty training sets, the loss function has no global minimum in \mathbb{R}^{p+1} . If there is a dividing hyperplane (=lin. separable), we can scale its parameters and increase the value of the loss function. Dies does not need to be the best dividing hyperplane! Fixed by regularizer, e.g. adding $\lambda \|w\|^2 + b^2$

3.5. Alternative Approach (Statistics)

Linear Model: $a = w_0 + w_1 x_1 + \dots + w_n x_n = w^T x$
 Probability: $\sigma(a) = \frac{1}{1 + e^{-a}}$ $D = \{(x_i, z_i)\}$
 Find $w_{MLE} = \text{arg max}_w P(D|w)$ $x_i \in \mathbb{R}^d, z_i \in \{0, 1\}$

$$P(D|w) = \prod_{i=1}^n \sigma(w^T x_i)^{z_i} (1 - \sigma(w^T x_i))^{1-z_i}$$

$L(w) = -\log(P(D|w))$, $\nabla_w L(w) = X(\sigma(X^T w) - z)$
 (Requires Hessian Matrix, which is positive semi-definite)
 Results are identical up to the factor $1/n$

4. Kernels

Kernels replace the standard inner products with some function that is a *suitable generalization of an inner product* to allow for nonlinear behavior.

A positive semidefinite Kernel is a function $\kappa: \mathbb{R}^{p \times} \rightarrow \mathbb{R}$ so that for all sets $\mathbf{X} = \{x_1, \dots, x_n\}$ the *Gram-Matrix* $K_{ij} = \kappa(x_i, x_j)$ is symmetric and positive semi definite

- Symmetric: $\kappa(x_i, x_j) = \kappa(x_j, x_i) \forall i, j$
- positive semi definite: $x^T K x \geq 0 \forall x$ or all eigenvalues ≥ 0

4.1. Common Kernels and Rules

- Linear kernel: $\kappa(x, y) = x^T y + c, c \geq 0$
- Polynomial Kernel $\kappa(x, y) = (ax^T y + c)^d, a, c, d \geq 0$
- Gaussian Kernel $\kappa(x, y) = \exp(-\frac{\|x-y\|^2}{2\sigma^2})$
- If κ_1, κ_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

5. Principal Component Analysis

Unsupervised Learning: Structure unlabeled data. Here: reduce dimension of input without losing too much information

Goal: Find $U_k \in \mathbb{R}^{p \times k}$ that minimizes $\sum_{i=1}^n \|x_i - U_k U_k^T x_i\|^2$

5.1. PCA - Singular Value Decomposition

Important: X : centered observation matrix SVD returns: $X = X D V^T$ with singular values $d_1 \geq d_2 \geq \dots \geq d_n$ on the diagonal of D

U_k satisfies goal with diagonal covariance matrix $S = U_k^T X X^T U_k$
 s_{ij} : j -th score of the i -th principal component. S : score matrix, U : loadings matrix

Cheaper than using $U_k^T X$ (nk operations instead $(2p-1)nk$)

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 $X_a X_a^T$ are the loadings of X_a only if X_a is centered

5.2. Statistical Interpretation

Goal: Find $Y = U^T X$ with $U^T U = I$ and declining variance of the components

5.3. Autoencoders

Autoencoder: Form of neural networks, where input is first mapped to a lower dimension k by f and then back to input $g \circ f(x_i) \approx x_i$. If f and g are linear and represented by $V \in \mathbb{R}^{k \times p}$ and $W \in \mathbb{R}^{p \times k}$ then $g \circ f(x_i) = W V x_i$ and the error is measured by the sum of squared distances $\sum_{i=1}^n \|x_i - X V x_i\|^2$, then the first k singular vectors of X are optimal.

WV has at most rank k

If f, g not linear: only approximation possible

6. Kernel-PCA

6.1. Linear PCA by inner products

If $K = X^T X$ the inner product matrix and $X = U^T$ the SVD of X , we can write $K = V \Sigma^T \Sigma V^T$ with $\Sigma^T \Sigma$ being diagonal. This is the eigenvalue decomposition of X

Let V_k and $\Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k)$ the first k eigenvectors and eigenvalues

Therefore: $U_k^T = \Sigma_k^{-1} V_k^T X^T$ 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 $x_n^T y = k_y$

If X is not centered, we can center K without centering X :

$$\tilde{K} = HKH \text{ with } H = (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T)$$

For new samples, we have to replace $\tilde{y} = y - \frac{1}{n} X \mathbf{1}_n$ and $\tilde{X} = XH$

$$U_k^T(\tilde{y}) = \Sigma_k^{-1} V_k^T \tilde{k}_y, \quad \tilde{k}_y = H k_y - \frac{1}{n} H K \mathbf{1}_n$$

6.2. Transition to Kernel PCA

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

$$k^{new} = [kappa(x_1, y) \dots kappa(x_n, y)]^T$$

$$k_{cent}^{new} = H k_{cent}^{new} - \frac{1}{n} H K \mathbf{1}_n$$

Kernel Principal Component Analysis

- for training set $\mathbf{X} = [x_1 \dots x_n]$, $x_i \in \mathbb{R}^p$
- Find suitable Kernel function $\kappa(\cdot)$ and compute Gram Matrix \mathbf{K}
- Compute centered Gram Matrix: $\tilde{K} = HKH$ with $H = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$
- Compute Eigenvalue Decomposition: $\tilde{K} = V \Lambda V^T$. Because \mathbf{K} is positive semi-definite and therefore the diagonal entries of Λ 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:
 $s_{new} = \Sigma_k^{-1} V_k^T k_{cent}^{new}$
 $k_{cent}^{new} = H k_{cent}^{new} - \frac{1}{n} H K \mathbf{1}_n$
 $k^{new} = [kappa(x_1, y) \dots kappa(x_n, y)]^T$

7. Feedforward Neural networks

Minimize any expected loss at the price of training many parameters
 Find $f \in \mathcal{F}$ that minimizes loss. If all functions in \mathcal{F} can be described by a set of k parameters $\Theta \in \mathbb{R}^k$, then the goal is for N samples:

$$\hat{\Theta} = \text{arg min}_{\Theta \in \mathbb{R}^k} \frac{1}{N} \sum_i L(f_{\Theta}(x_i))$$

One layer of a FNN consists of a linear function $\phi_W(h) = W h$ and a nonlinear *activation function*, e.g. the ReLU: $\sigma(t) = \max\{0, t\}$:

$$f: \mathbb{R}^p \rightarrow \mathbb{R}^k: \sigma_1 \circ \phi_{W_1} \circ \dots \circ \sigma_l \circ \phi_{W_l}(X)$$

With p/k Input/Output dimension and l layers, *Deep FNN*: $l > 3$
 k depends on loss function (Output of FNN is input of Loss)

7.1. Training FNNs (Backpropagation)

Functions $g: \mathbb{R}^k \rightarrow \mathbb{R}^l$ and $h: \mathbb{R}^l \rightarrow \mathbb{R}^m$ with g being diff'able at x and h at $y = g(x)$ with Jacobi matrices $J_g(x)$ and $J_h(y)$ then
 $h \circ g: \mathbb{R}^k \rightarrow \mathbb{R}^m$ is diff'able at x with
 $J_{h \circ g}(x) = J_h(g(x)) \cdot J_g(x)$

8. Support Vector Machine (Basic Linear SVM)

- Idea: Find hyperplane to divide x, y into two subspaces
- Hyperplane: $\mathcal{H}_{w,b} = \{x \in \mathbb{R}^d | w^T x - b = 0\}$, w normal to \mathcal{H}
- Euclidean distance from x to $\mathcal{H}_{w,b}$: $\delta(x, \mathcal{H}_{w,b}) = \frac{w^T x - b}{\|w\|}$
- Goal: Find $y_i(w^T x_i - b) \geq 1 \forall i, 1, \dots, N$
- Optimization Problem: $\min \frac{1}{2} \|w\|^2$ that fulfills Goal

8.1. Karush Kuhn Tucker Conditions and Linear SVMs

$$\min_{w, b, \lambda \geq 0} L(w, b, \lambda)$$

$$L(w, b, \lambda) = \frac{1}{2} \|w\|^2 - \sum_i \lambda_i y_i (w^T x_i - b) + \lambda_i$$

$$\nabla_{(w,b)} L(w, b, \lambda) = \begin{bmatrix} w - \sum_i \lambda_i y_i x_i \\ -\sum_i \lambda_i y_i \\ w^* = \sum_i \lambda_i y_i x_i \quad \sum_i \lambda_i y_i = 0 \end{bmatrix}$$

$$L_D(\lambda) = \sum_i \lambda_i - \frac{1}{2} \lambda^T H \lambda \quad \text{With } h_{ij} = y_i y_j x_i^T x_j$$

New convex quadratic opt problem: $\lambda^* = \max_{\lambda} L_D(\lambda)$

- 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 LR
- SVM is preferred when using kernels compared to LR

9. Decision Making

- **Decision making** is about mapping an input vector to an output value.
- **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 .
- **Loss** is a positive, real measure of the deviation between a target and an output from a decision function.
- **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 global decision function f out of class F that minimises the expected prediction error.
- The aim of **local methods** is to find the best output which minimizes the expectation of loss on the distribution of Y conditioned on known X .
- We can obtain 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, and so we take the **k -nearest neighbours** to x as approximations of x , thus giving us a set of samples representing the distribution of values $X=x$ such that local method for finding output can now be applied to these samples.

10. 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.
- As the number of dimensions increases, the number of samples required to obtain an accurate representation of the probability distribution also increases a lot.
- The **empty space theorem** says that as dimensionality increases, samples tend to be located at the tail end of their distributions.
- As the number of dimensions increases, the 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 β away from the center: (For large p this becomes 1)

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

11. 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, g convex, then: $\max(f, g), f + g, g \circ h$ (if g non-decreasing) are convex. Local minimum of *strictly* convex is global minimum and unique

12. Logistic Regression

12.1. 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 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ is also convex, and can therefore be easily optimized via minimization.

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

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

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

Gradient-based methods can be used to find affine weights and bias that optimize 3.3.

12.2. Overfitting on Logistic Regression

- When the samples are linearly separable, 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:
 $y_i(\mathbf{w}_s^T \mathbf{x}_i + b_s) < 0 \forall i$
 then the lossfunction has no global minimum in \mathbb{R}^{p+1}

12.2.1. Proof

This is proven by noting that for a global minimum to exist, a point with small, positive loss exists (3.17) such that all points have losses greater than or equal to this point (3.16). This has to be satisfied along the linearly separable constraint (3.18). If this constraint is always satisfied, then (3.19) is always positive, which when subbed into (3.15) and scaling weights and bias to infinity we find that the loss approaches zero (3.21), thus violating the global minimum condition (3.18). This is prevented via regularization (3.23).

13. Principal Component Analysis

PCA represents the data projected onto a lower-dimensional plane which captures most of the variance of the dataset.

13.1. Geometric Interpretation

- Imagining that \mathbf{x} is a vector in p dimensions, and U_k describes a lower dimensional subspace, then the coordinates of \mathbf{x} projected onto U_k is found by first obtaining the scores of \mathbf{x} along U_k and then taking the linear combination of these scores and the U_k vectors: $\pi_U(\mathbf{x}) = U_k U_k^T \mathbf{x}$, $U_k \in \mathbb{R}^{p \times k}$
- 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.

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

- The dataset can be represented as the dot product of an orthonormal matrix, a diagonal matrix and another orthonormal matrix via singular value decomposition: $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$.

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).

13.2. 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)

13.3. 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). Thus we can push U inside the covariance formula (4.8) to show that the scores with uncorrelated dimensions are obtained by projecting X onto U . (4.9)

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

13.4. Error Model Interpretation

- The error model interpretation is that the samples \mathbf{X} is obtained via the addition between a signal matrix L and a noise matrix N , where L is lying on a lower dimensional subspace than \mathbf{X} , and the goal is to try and find L . Theorem 4.2
- Setting L equal to the SVD (p dimensions with k) minimizes the Frobenius norm of the difference between \mathbf{X} and L given that the dimension k of the subspace which L lies in is known beforehand.

13.5. Relation to Autoencoders

- Autoencoders map an input to lower 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 $g \circ f(\mathbf{x}_i) \approx \mathbf{x}_i$
- Letting f and g be linear mappings represented by matrices, the reconstruction error $J(\mathbf{W}, \mathbf{V}) = \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{W} \mathbf{V} \mathbf{x}_i\|_2^2$ 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 \mathbf{X} , the $\mathbf{V} = U_k^T$ and $\mathbf{W} = U_k$ 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.

14. Feed forward Neural Networks

14.1. 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)

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

14.2. 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).

$\mathbf{J}_{\mathbf{W} \circ \mathbf{g}}(\mathbf{x}) = \mathbf{W} \cdot \mathbf{J}_{\mathbf{g}}(\mathbf{x})$ (5.6)
 $\mathbf{J}_{(j), \sigma}(\mathbf{x}) = \begin{cases} \sigma'(x_i) & \text{if } i = j \\ 0 & \text{else} \end{cases}$ (5.7)
 $\mathbf{J}_{mult}(\mathbf{x}) = \text{diag}(\mathbf{x}^T) \in \mathbb{R}^{m \times mn}$ (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):
 $\frac{\partial}{\partial \mathbf{W}_L} F = \mathbf{J}_L(\mathbf{h}_L) \cdot \mathbf{J}_{\sigma_L}(\mathbf{W}_L \mathbf{h}_{L-1}) \cdot \mathbf{J}_{mult}(\mathbf{h}_{L-1})$
 $\frac{\partial}{\partial \mathbf{W}_j} F = \mathbf{J}_L(\mathbf{h}_L) \cdot \mathbf{J}_{\sigma_L}(\mathbf{W}_L \mathbf{h}_{L-1}) \cdot \dots \cdot \mathbf{J}_{\sigma_j}(\mathbf{W}_j \mathbf{h}_{j-1}) \cdot \mathbf{J}_{mult}(\mathbf{h}_{j-1})$ for the j -th layer

15. Kernel Trick

- A kernel is a function that maps two sets containing real numbers (representing a sample) to a single scalar, such that for all possible sets of samples (contained in a data matrix with columns corresponding to samples) the Gram-Matrix (6.1) is symmetric and positive semidefinite.

A positive semidefinite Kernel is a function $\kappa : \mathbb{R}^{p \times p} \rightarrow \mathbb{R}$ so that for all 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 semi definite

- 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 psd (since psd eigenvalues are nonnegative and the determinant is the product of the eigenvalues), but we cannot say anything if the determinant is positive.

16. Kernel PCA

16.1. 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).
- 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): $\mathbf{K} = \mathbf{V}\Sigma^T\mathbf{V}^T$
Centered Gram Matrix (7.8): $\tilde{\mathbf{K}} = \mathbf{H}\mathbf{K}\mathbf{H}$ with

$$\mathbf{H} = (\mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^T)$$

Computing scores for new sample \mathbf{y} : $\mathbf{U}_k^T(\tilde{\mathbf{y}}) = \Sigma_k^{-1}\mathbf{V}_k^T\tilde{\mathbf{k}}_y$.

$$\tilde{\mathbf{k}}_y = \mathbf{H}\mathbf{k}_y - \frac{1}{n}\mathbf{H}\mathbf{K}\mathbf{1}_n$$

16.2. 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)

17. Support Vector Machines

17.1. Geometry

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.

$$\text{Hyperplane: } \mathcal{H}_{w,b} = \{x \in \mathbb{R} \mid w^T x - b = 0\} \quad (8.1)$$

$$w \text{ is normal to } \mathcal{H}: w^T(x_1 - x_2) = 0 \quad (8.2)$$

$$\text{Signed Distance: } \delta(x, \mathcal{H}_{w,b}) = \frac{w^T x - b}{\|w\|} \quad (8.3)$$

$$\text{Margin: } \mathcal{H}_+ = \{x \in \mathbb{R} \mid w^T x - b = 1\}$$

$$\mathcal{H}_- = \{x \in \mathbb{R} \mid w^T x - b = -1\} \quad (8.8-9)$$

- 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).

17.2. 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).

$$\text{Conditions: } w^T x_i - b \geq +1 \text{ for } y_i = 1 \text{ and}$$

$$w^T x_i - b \geq -1 \text{ for } y_i = -1$$

$$\text{Compact: } y_i \cdot (w^T x_i - b) \geq 1 \quad \forall i \quad (8.12) \text{ Minimization:}$$

$$\min_w \frac{1}{2} \|w\|_2^2 \text{ s.t. } y_i (w^T x_i - b) \geq 1 \quad \forall i$$

17.3. KKT and Lagrangian Duality

The Karush-Kuhn Tucker conditions says 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):

$$\min_{z \in \mathbb{R}^n} f(z) \text{ s.t. } c_i(z) = 0 \quad \forall i \in \mathcal{E} \text{ and s.t. } c_j(z) \geq 0 \quad \forall j \in \mathcal{I}$$

$$\text{Lagrange Function (8.16): } L(z, \lambda) = f(z) - \sum_{i \in \mathcal{I} \cup \mathcal{E}} \lambda_i c_i(z)$$

Karush-Kuhn-Tucker Conditions

For z^* as a solution to 8.15, there exists a Lagrange multiplier λ^* such that:

$$\begin{aligned} -\Delta_z L(z^*, \lambda^*) &= 0 \\ -c_i(z^*) &= 0 \quad \forall i \in \mathcal{E} \\ -c_j(z^*) &\geq 0 \quad \forall j \in \mathcal{I} \\ -\lambda_i^* &\geq 0 \quad \forall i \in \mathcal{I} \\ -\lambda_i^* c_i(z^*) &= 0 \quad \forall i \in \mathcal{I} \cup \mathcal{E} \end{aligned}$$

- 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).

17.4. 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 (8.26 to 8.29).

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

- Then the dual function needs to be calculated by taking the infimum of 8.25, which is accomplished by substituting the gradients w.r.t w and b (8.26) inside, resulting in the dual function (8.30a).
- Maximising this function w.r.t. its constraints is then the dual problem for SVM. (8.30)

18. Useful Facts

The matrix resulting from the dot product of between an $R^{n \times p}$ matrix and $R^{p \times n}$ has at most a rank of p .

19. Homework and Assignments

Given a prediction table, eg:

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

Calculate Covariance Matrix from Table

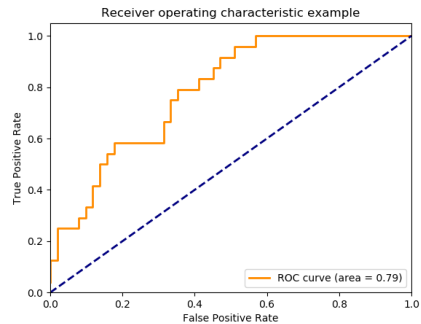
- Calculate Means for X_1 and X_2 : μ_1, μ_2
- Create X-Matrix, e.g.: $\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$
- Create p-Matrix: $\begin{bmatrix} p_X(0,0) & 0 & 0 & 0 \\ 0 & p_X(1,0) & 0 & 0 \\ 0 & 0 & p_X(0,1) & 0 \\ 0 & 0 & 0 & p_X(1,1) \end{bmatrix}$
- Calculate Covariance: $Cov = Xp_X^T$

19.1. Classification Analysis

19.1.1. 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:

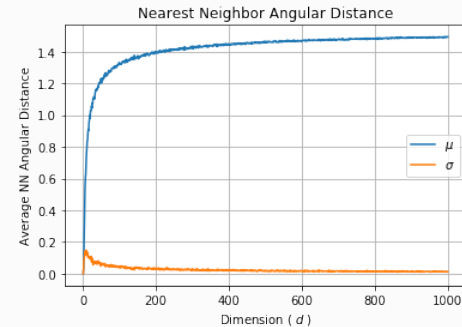
- The relationship between sensitivity and specificity. For example, a decrease in sensitivity results in an increase in specificity.
- Test accuracy; the closer the graph is to the top and left-hand borders, the more accurate the test. Likewise, the closer the graph to the diagonal, the less accurate 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.



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).

19.2. 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 space 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 .



19.3. 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

19.4. 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)

19.4.1. 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 ϵ , i.e.

$$\frac{\|\mathbf{X} - \mathbf{P}\mathbf{X}\|_F^2}{\|\mathbf{X}\|_F^2} = \frac{\sum_{i=1}^M \|\mathbf{x}_i - \mathbf{P}\mathbf{x}_i\|_2^2}{\sum_{i=1}^N \|\mathbf{x}_i\|_2^2} \leq \epsilon,$$

where ϵ 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)

