

DSA5205 Data Science for Quantitative Finance

AY2024/25 Sem1 By Zhao Peiduo

Introduction and Background

Big Data and "V"s

Big Data is characterized by several "V"s:

- **Volume:** Massive amounts of data
- **Velocity:** The speed of data processing
- **Variety:** Different types of data (structured/unstructured)
- **Veracity:** The quality or trustworthiness of the data
- **Value:** The potential benefit derived from analyzing the data

Machine Learning Overview

- **Supervised Learning:** Learn a function from labeled data
- **Unsupervised Learning:** Find patterns without labeled data
- **Reinforcement Learning:** Learn from rewards and punishments

Quantitative Finance

Quantitative finance uses statistical and mathematical models to analyze financial markets and manage risks. Common models include:

- **CAPM:** Capital Asset Pricing Model
- **Black-Scholes:** Option pricing model

The Data Science Process

1. Problem definition
2. Data collection
3. Data cleaning and preprocessing
4. Model building (Machine learning, statistics)
5. Evaluation and interpretation
6. Reporting and visualization

Challenges in Data Science

- **Data Cleaning:** Handling missing values, outliers
- **Model Selection:** Choosing the right model
- **Ambiguity:** Dealing with uncertainty in data

Distribution and Risks

Moments

Let X be a random variable. The k^{th} moment of X is defined as: $E(X^k)$

The first moment is the expectation. The k^{th} central moment is: $\mu_k = E[(X - E(X))^k]$

Variance:

Skewness

Skewness is a measure of symmetry. For a continuous random variable Y : $Sk(Y) = E\left[\frac{(Y - E(Y))^3}{\sigma^3}\right]$

For a discrete random variable: $Sk(Y) = \sum \frac{(y - E(Y))^3}{\sigma^3} f(y)$

Kurtosis

Kurtosis measures tail thickness. The kurtosis of a random variable Y is: $Kur(Y) = E\left[\frac{(Y - \mu_Y)^4}{\sigma_Y^4}\right]$

For a normal distribution $Y \sim N(\mu, \sigma^2)$, the kurtosis is 3. The excess kurtosis is: $Kur(Y) - 3$

Heavy-Tailed Distributions

A distribution is heavy-tailed if its tails are thicker than the normal distribution. A right Pareto tail has the form:

$f(y) \sim Ay^{-(a+1)}$ as $y \rightarrow \infty$

The parameter a is called the tail index.

Student's t-Distributions

The probability density function (pdf) of the t-distribution with ν degrees of freedom is: $f(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$

Mean: $E[X] = 0$ for $\nu > 1$

Variance: $Var(X) = \frac{\nu}{\nu-2}$ for $\nu > 2$

Quantile-Based Measures

Quantiles of a distribution: Let $X \sim F(x)$. The p^{th} quantile of $F(x)$ is defined as: $q_p = F^{-1}(p)$

Interquartile range (IQR): $IQR = q_{0.75} - q_{0.25}$

Risk Measures

Value at Risk (VaR) is defined as: $VaR_\alpha = q_\alpha(F_L)$ such that $P(L \geq VaR_\alpha) = \alpha$

Expected Shortfall (ES) is: $ES_\alpha = E(L|L \geq VaR_\alpha)$

Jarque-Bera Test

The Jarque-Bera test is a test of normality based on skewness and kurtosis. For a sample Y_1, \dots, Y_n :

$$Sk = \frac{1}{n} \sum_{i=1}^n \left(\frac{Y_i - \bar{Y}}{s} \right)^3$$
$$Kur = \frac{1}{n} \sum_{i=1}^n \left(\frac{Y_i - \bar{Y}}{s} \right)^4$$
$$JB = n \left(\frac{Sk^2}{6} + \frac{(Kur - 3)^2}{24} \right)$$

Under the null hypothesis of normality, the test statistic JB follows a chi-squared distribution with 2 degrees of freedom.

Heavy-Tailed and Skewed Distributions

A generalized error distribution (GED) with shape parameter ν is used for modeling heavy tails and skewness:

$$f_{\text{GED}}(y|\nu) = k(\nu) \exp\left(-\frac{1}{2} \left| \frac{y}{\lambda\nu} \right|^\nu\right), \quad -\infty < y < \infty$$

The tail behavior is controlled by ν , with smaller values indicating heavier tails.

Profile Likelihood

Profile likelihood is used for constructing confidence intervals by fixing one parameter at a time and maximizing over the others:

$$L_{\max}(\theta_1) = \max_{\theta_2} L(\theta_1, \theta_2)$$

Generalized Error Distributions (GED)

The Generalized Error Distribution (GED) has a flexible tail weight controlled by the shape parameter ν .

If $Y \sim GED(\mu, \sigma, \nu)$, its pdf is:

$$f_{\text{GED}}(y|\nu) = k(\nu) \exp\left(-\frac{1}{2} \left| \frac{y}{\lambda\nu} \right|^\nu\right)$$

The distribution is symmetric about μ , with $\nu = 2$ corresponding to the normal distribution, and $\nu = 1$ being the double-exponential. The tail weight increases as ν decreases.

Quantile-Quantile (QQ) Plot

QQ plots are used to compare the quantiles of a dataset to those of a theoretical distribution, typically the normal distribution.

If the data follows the theoretical distribution, the points should lie approximately on the 45-degree line. Deviations indicate skewness, heavy tails, or other non-normal behavior.

Fisher Information

Fisher information is defined as the expected value of the negative second derivative of the log-likelihood function.

For a parameter θ :

$$I(\theta) = -E\left[\frac{\partial^2}{\partial \theta^2} \log L(\theta)\right]$$

Fisher information quantifies the amount of information that an observable random variable Y carries about an unknown parameter θ . The standard error of an estimator is inversely proportional to the square root of the Fisher information.

Central Limit Theorem for the MLE

The CLT for the Maximum Likelihood Estimator (MLE) states:

$$\hat{\theta} \sim N\left(\theta, I(\hat{\theta})^{-1}\right)$$

Large-sample confidence interval for the MLE of θ is given by:

$$\hat{\theta} \pm s_{\hat{\theta}} z_{\alpha/2}$$

Use the observed Fisher information to approximate $I(\theta)$:

$$I^{obs}(\theta) = -\left[\frac{d^2}{d\theta^2} \log L(\theta)\right]$$

The standard error is given by:

$$s_{\hat{\theta}}^{obs} = \frac{1}{\sqrt{I^{obs}(\hat{\theta})}}$$

Replace $s_{\hat{\theta}}^A$ with $s_{\hat{\theta}}^{obs}$ in the confidence interval for more convenience and accuracy. Fisher information matrix for the multivariate case is replaced by the Hessian matrix.

MLE has a small bias:

$$BIAS(\hat{\theta}_{ML}) = E(\hat{\theta}_{ML}) - \theta \sim -\frac{A}{n}, \quad \text{as } n \rightarrow \infty$$

Variance of the MLE:

$$VAR(\hat{\theta}_{ML}) \sim \frac{B}{n}, \quad \text{as } n \rightarrow \infty$$

Likelihood Ratio Tests (LRT)

Suppose that θ is a parameter vector,

$$H_0 : \theta = \theta_0 \quad \text{vs} \quad H_1 : \theta \neq \theta_0$$

The rejection region is the set of possible samples that lead us to reject H_0 . We need to keep the probability of a type I error below a prespecified small value called the significance level α .

Let $\hat{\theta}_{ML}$ be the MLE without restrictions (full model) and $\hat{\theta}_{0,ML}$ be the value of θ that maximizes $L(\theta)$ under H_0 (reduced model). The likelihood ratio test rejects H_0 if:

$$2 \left[\log L(\hat{\theta}_{ML}) - \log L(\hat{\theta}_{0,ML}) \right] \geq c$$

where c is a critical value which gives a level that is exactly equal to α or $c = \chi_{\alpha, m}^2$, the α -upper quantile value of the chi-squared distribution with m degrees of freedom.

Robust Estimation

Robust estimators are resistant to outliers and deviations from model assumptions.

An example is the trimmed mean, where a proportion of the largest and smallest values are excluded. Another is the Median Absolute Deviation (MAD), which is a robust measure of scale:

$$\hat{\sigma}_{MAD} = 1.4826 \times \text{median}(|Y_i - \text{median}(Y)|)$$

MAD is less sensitive to extreme values than standard deviation.

Parsimonious

Tradeoff between bias and variance. The optimal situation is little bias without excessive parameters.

AIC and BIC

The Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) are used to compare models. Both criteria balance goodness of fit with model complexity:

$$AIC = -2 \log L(\hat{\theta}_{MLE}) + 2p$$

$$BIC = -2 \log L(\hat{\theta}_{MLE}) + p \log n$$

Where p is the number of parameters and n is the sample size. Lower values of AIC or BIC indicate better models, but BIC penalizes complexity more than AIC.

Variance-Covariance Method

In the variance-covariance method, the loss distribution is assumed to be multivariate normal. The linearized loss is approximated by:

$$L \sim N(\mu, \sigma^2)$$

Value at Risk (VaR) and Expected Shortfall (ES) are estimated from the distribution. For normal distribution:

$$VaR_{\alpha} = \mu + \sigma \Phi^{-1}(1 - \alpha)$$

$$ES_{\alpha} = \mu + \sigma \frac{\phi(\Phi^{-1}(1 - \alpha))}{\alpha}$$

where ϕ is the standard normal density and Φ^{-1} is the inverse cumulative distribution function (quantile) of the normal distribution.

Hill Estimator

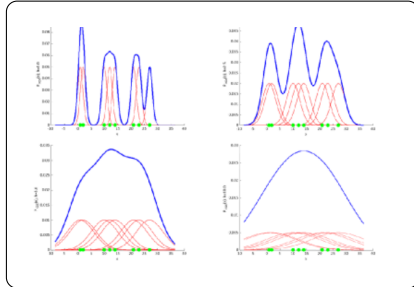
The Hill estimator is used to estimate the tail index of heavy-tailed distributions. For order statistics $X_{(1)}, X_{(2)}, \dots, X_{(n)}$, the Hill estimator is defined as:

$$\hat{\gamma}_H = \frac{1}{k} \sum_{i=1}^k \log \left(\frac{X_{(i)}}{X_{(k+1)}} \right)$$

The parameter k determines how many of the extreme values (largest observations) are used. A Hill plot helps to determine the appropriate value of k . The Hill estimator is particularly useful for Pareto-like heavy-tailed distributions.

Bandwidth Selection for Smoothing Kernels:

- **High bandwidth (h):** Leads to *over-smoothing*, which can mask the data distribution.
- **Low bandwidth (h):** Results in spiky and *hard-to-interpret* density estimation.



Nonparametric Methods and Resampling

Histogram

A histogram divides the sample space into bins and approximates the density at each bin's center:

$$p_H(X) = \frac{\text{Number of } x(k) \text{ in the same bin as } x}{\text{Width of bin}}$$

Hyperparameters: bin width and starting position of the first bin.

Kernel Density Estimation (KDE)

KDE estimates the probability density of a random variable, where $K(u)$ is the kernel function, and h is the bandwidth (smoothing parameter):

$$\hat{p}_{KDE}(x) = \frac{1}{nhD} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

Non-Parametric Density Estimation

The general expression for non-parametric density estimation becomes:

$$p(x) = \frac{k}{NV}$$

V volume surrounding x
 N total number of examples
 k number of examples inside V

Parzen Windows

For Parzen window estimation, the kernel function is:

$$K(u) = \begin{cases} 1 & \text{if } |u| \leq \frac{1}{2}, \forall j = 1 \dots D \\ 0 & \text{otherwise.} \end{cases}$$

To understand the role of the kernel function, we compute the expectation of the estimate $p_{KDE}(X)$:

$$\mathbb{E}[p_{KDE}(x)] = \frac{1}{NhD} \sum_{n=1}^N \mathbb{E} \left[K \left(\frac{x - x^{(n)}}{h} \right) \right] = \frac{1}{hD} \mathbb{E} \left[K \left(\frac{x - x^{(n)}}{h} \right) \right] = \frac{1}{hD} \int K \left(\frac{x - x'}{h} \right) p(x') dx'.$$

- Where we have assumed that vectors $x^{(n)}$ are drawn independently from the true density $p(x)$.
- We can see that the expectation of $p_{KDE}(X)$ is a convolution of the true density $p(x)$ with the kernel function.
 - Thus, the kernel width h plays the role of a smoothing parameter: the wider h is, the smoother the estimate $p_{KDE}(X)$.
 - For $h \rightarrow 0$, the kernel approaches a Dirac delta function, and $p_{KDE}(X)$ approaches the true density.
 - However, in practice, we have a finite number of observations, so h cannot be made arbitrarily small since the density estimate $p_{KDE}(X)$ would then degenerate to a set of impulses located at the training data points.

Smooth Kernels

The Parzen window yields density estimates that have discontinuities and weights equally all points x_i , regardless of their distance to the estimation point. For these reasons, the Parzen window is commonly replaced with a smooth kernel function $K(u)$

Let $K(u)$ be Unimodal pdf, such as the Gaussian, then its corresponding density estimate is:

$$K(x) = (2\pi)^{-D/2} e^{-1/2x'^2}$$
$$p_{KDE}(x) = \frac{1}{NhD} \sum_{n=1}^N K \left(\frac{x - x^{(k)}}{h} \right)$$

Usually, but not always, $K(u)$ will be radially symmetric and $\int_{\mathbb{R}^D} K(x) dx = 1$

Bandwidth Selection

The optimal bandwidth h minimizes the mean squared error (MSE) between the KDE and the true density:

$$MSE = \mathbb{E}[(\hat{p}_{KDE}(x) - p(x))^2] = \text{Bias}^2 + \text{Variance}.$$

Assume a standard density function and find the value of the bandwidth that minimizes the integral of the square error (MISE):

$$h_{MISE} = \arg \min \left\{ \mathbb{E} \left[\int (p_{KDE}(x) - p(x))^2 dx \right] \right\}.$$

The optimal bandwidth for a Gaussian kernel is:

$$h^* = 1.06 \cdot \sigma \cdot n^{-1/5}$$

Or we can use the more robust inter quantile range (IQR) instead of sample variance:

$$h^* = 0.9AN^{-1/5} \quad \text{where} \quad A = \min \left(\sigma, \frac{\text{IQR}}{1.34} \right)$$

Maximum Likelihood Cross-validation

- The ML estimate of h is degenerate since it yields $h_{ML} = 0$, a density estimate with Dirac delta functions at each training data point
- A practical alternative is to maximize the "pseudo-likelihood" computed using leave-one-out cross-validation

$$h^* = \arg \max \left\{ \frac{1}{N} \sum_{n=1}^N \log p_{-n} \left(x^{(n)} \right) \right\}$$

Where

$$p_{-n} \left(x^{(n)} \right) = \frac{1}{(N-1)h} \sum_{m=1, m \neq n}^N K \left(\frac{x^{(n)} - x^{(m)}}{h} \right)$$

Multivariate Density Estimation

$$p_{KDE}(x) = \frac{1}{Nh^D} \sum_{n=1}^N K \left(\frac{x - x^{(n)}}{h} \right)$$

- Notice that the bandwidth h is the same for all the axes, so this density estimate will weight all the axes equally.
- If one or several of the features has larger spread than the others, we should use a vector of smoothing parameters or even a full covariance matrix, which complicates the procedure.

There are two basic alternatives to solve the scaling problem without having to use a more general KDE:

- Pre-scaling each axis (normalize to unit variance, for instance).
- Pre-whitening the data (linearly transform so $\Sigma = I$), estimate the density, and then transform back [Fukunaga].
 - The whitening transform is $y = \Lambda^{-1/2} M^T x$, where Λ and M are the eigenvalue and eigenvector matrices of Σ .
 - Fukunaga's method is equivalent to using a hyper-ellipsoidal kernel.

Product Kernels

$$p_{PKDE}(x) = \frac{1}{N} \sum_{i=1}^N K(x, x^{(n)}, h_1, \dots, h_D)$$

where

$$K(x, x^{(n)}, h_1, \dots, h_D) = \frac{1}{h_1 \dots h_D} \prod_{d=1}^D K_d \left(\frac{x_d - x_d^{(n)}}{h_d} \right)$$

- The product kernel consists of the product of one-dimensional kernels:
 - Typically, the same kernel function is used in each dimension ($K_d(x) = K(x)$), and only the bandwidths are allowed to differ.
 - Bandwidth selection can then be performed with any of the methods presented for univariate density estimation.
- Note that although $K(x, x^{(n)}, h_1, \dots, h_D)$ uses kernel independence, this does not imply we assume the features are independent:
 - If we assumed feature independence, the density estimation (DE) would have the expression:

$$p_{\text{FEAT-IND}}(x) = \prod_{d=1}^D \frac{1}{Nh_d} \sum_{i=1}^N K_d \left(\frac{x_d - x_d^{(i)}}{h_d} \right).$$

Transformation KDE

Algorithm

1. Start with data Y_1, \dots, Y_n .
2. Transform the data to $X_1 = g(Y_1), \dots, X_n = g(Y_n)$; g is monotonic.
3. Let f_X be the usual KDE using X_1, \dots, X_n .
4. $f_Y(y) = f_X \{g(y)\} \left| g'(y) \right|$. Plugging in $y = g^{-1}(x)$, then

$$f_Y(g^{-1}(x)) = f_X \{x\} \left| g'(g^{-1}(x)) \right|$$

K-fold Cross-validation in Detail

Let the K parts be C_1, C_2, \dots, C_K , where C_k denotes the indices of the observations in part k . There are n_k observations in part k : if n is a multiple of K , then $n_k = n/K$.

$$\text{Compute } CV_{(K)} = \sum_{k=1}^K \frac{n_k}{n} MSE_k$$

where $MSE_k = \sum_{i \in C_k} (y_i - \hat{y}_i)^2 / n_k$ and \hat{y}_i is the fit for observation i , obtained from the data with part k removed.

Setting $K = n$ yields n -fold or *leave-one-out cross-validation* (LOOCV).

With least-squares linear or polynomial regression, an amazing shortcut makes the cost of LOOCV the same as that of a single model fit! The following formula holds:

$$CV_{(K)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - h_i} \right)^2$$

where \hat{y}_i is the i th fitted value from the original least squares fit, and h_i is the leverage (diagonal of the "hat" matrix; see ESL book for details). This is like the ordinary MSE, except the i th residual is divided by $1 - h_i$.

Cross-Validation

K -fold cross-validation divides the data into K parts. The model is trained on $K - 1$ parts and tested on the k -th part:

$$CV(K) = \frac{1}{n} \sum_{k=1}^K \frac{1}{n_k} \sum_{i \in C_k} (\hat{y}_i - y_i)^2.$$

Bootstrap

Bootstrap estimates the uncertainty of an estimator by resampling:

$$SE_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^B (\hat{\alpha}_r^* - \hat{\alpha}^*)^2}.$$

Multivariate Methods and Factor Models

Multivariate Data and Factor Models

Covariance Matrices

Covariance measures the direction but not the strength of a linear relationship between two random variables X and Y . The covariance matrix of a random vector $\mathbf{Y} = (Y_1, \dots, Y_d)$ is defined as:

$$\text{COV}(\mathbf{Y}) = \mathbb{E} \left[(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^\top \right].$$

The diagonal elements are the variances of individual components.

Multivariate Normal Distribution

A random vector $\mathbf{Y} = (Y_1, \dots, Y_d)$ follows a multivariate normal distribution if its probability density function (PDF) is:

$$\phi_d(\mathbf{y}; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right).$$

Here, $\boldsymbol{\mu}$ is the mean vector and Σ is the covariance matrix.

Principal Component Analysis (PCA)

PCA reduces the dimensionality of multivariate data by finding new orthogonal axes (principal components) that capture the maximum variance. The principal components are eigenvectors of the covariance matrix, and the corresponding eigenvalues represent the variance explained by each component.

Multivariate t-Distribution

The random vector \mathbf{Y} follows a multivariate t-distribution if:

$$\mathbf{Y} = \boldsymbol{\mu} + \sqrt{\frac{\nu}{W}} \mathbf{Z},$$

where W follows a chi-squared distribution with ν degrees of freedom, and \mathbf{Z} follows a multivariate normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix Σ .

Copulas

A copula is a multivariate distribution function where the marginal distributions are uniform. The copula function links the marginal distributions to form a joint distribution. For a random vector $\mathbf{Y} = (Y_1, \dots, Y_d)$, the copula is defined as:

$$C_{\mathbf{Y}}(u_1, \dots, u_d) = \mathbb{P}(F_{Y_1}(Y_1) \leq u_1, \dots, F_{Y_d}(Y_d) \leq u_d),$$

where $F_{Y_i}(Y_i)$ are the marginal cumulative distribution functions (CDFs) of Y_i .

Linear Regression

Linear regression models the relationship between the outcome Y and predictors X_1, X_2, \dots, X_d :

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_d X_d + \epsilon$$

The coefficients β are estimated by minimizing the sum of squared residuals (errors):

$$J(\beta) = \frac{1}{2n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Gradient Descent

To optimize the cost function $J(\beta)$, gradient descent iteratively updates the parameters β as follows:

$$\beta_j \leftarrow \beta_j - \alpha \frac{\partial J}{\partial \beta_j}$$

Here, α is the learning rate. Gradient updates:

$$\frac{\partial J}{\partial \beta_j} = \frac{1}{n} \sum_{i=1}^n (h_{\beta}(X^{(i)}) - Y^{(i)}) X_j^{(i)}$$

Polynomial Regression

Extends linear regression by allowing higher-order terms:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_p X^p + \epsilon$$

This increases flexibility but can lead to overfitting.

Generalized Additive Models (GAM)

GAMs are an extension of linear models where each predictor X_j is modeled with a smooth function:

$$Y = \beta_0 + f_1(X_1) + f_2(X_2) + \dots + f_p(X_p) + \epsilon$$

GAMs can capture nonlinear relationships without specifying the exact form.

Model Selection Overview

Subset Selection

Subset selection identifies a subset of predictors that best relate to the response. The best subset is chosen based on criteria like:

- Residual Sum of Squares (RSS)
- R^2
- AIC, BIC, Cp, or cross-validation error

Best Subset Selection:

1. Start with the null model M_0 which has no predictors.
2. For each k , fit models with exactly k predictors and select the one with the smallest RSS or highest R^2 .
3. Use cross-validation or another criterion to choose the best model among all.

Stepwise Selection

Stepwise selection simplifies the search process:

- **Forward Stepwise:** Start with no predictors and iteratively add the one that improves the model the most.
- **Backward Stepwise:** Start with all predictors and iteratively remove the least useful.

Regularization and Shrinkage Methods

Shrinkage methods penalize the model's complexity, reducing variance:

- **Ridge Regression:** Shrinks coefficients by adding a penalty proportional to their squared values:

$$\hat{\beta} = \arg \min \left[RSS + \lambda \sum_{j=1}^p \beta_j^2 \right]$$

- **Lasso:** Uses an L_1 penalty to encourage sparsity, setting some coefficients to zero:

$$\hat{\beta} = \arg \min \left[RSS + \lambda \sum_{j=1}^p |\beta_j| \right]$$

Choosing the Optimal Model

Model selection is based on minimizing test error, often estimated via cross-validation or information criteria like AIC and BIC:

$$AIC = -2 \log(L) + 2 \cdot d \quad \text{and} \quad BIC = -2 \log(L) + \log(n) \cdot d$$

Here, L is the likelihood of the model, and d is the number of parameters.

Cross-Validation

- Split the data into K folds.
- Train the model on $K - 1$ folds and validate on the remaining fold.
- Repeat for each fold and average the validation errors.

This procedure provides a direct estimate of test error.

Time Series Basics

- A **time series** is sequential data indexed over time.
- A **stochastic process** is a theoretical construct for time series, where data points are random variables.

Stationarity

- **Stationarity:**
 - Time-invariant statistical properties.
 - A process is stationary if $P(Y_t = y) = P(Y_{t+k} = y)$.
- **Stationary Process:**
 - Oscillates around a mean (mean-reversion).
 - Example: AR(1) stationarity condition: $-1 < \phi < 1$.

Autocorrelation

- **Autocorrelation Function (ACF):** Measures linear dependence between observations k -time units apart.
- **Autocovariance:**

$$\gamma_k = \mathbb{E}[(Y_t - \mu)(Y_{t-k} - \mu)]$$

AR and MA Models

- **AR(p)** (Autoregressive Model):

$$Y_t = \delta + \phi_1 Y_{t-1} + \dots + \phi_p Y_{t-p} + \epsilon_t$$

- **MA(q)** (Moving Average Model):

$$Y_t = \mu + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$$

- **ARMA(p,q)** (Combination of AR and MA):

$$Y_t = \delta + \phi_1 Y_{t-1} + \dots + \phi_p Y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$$

Trend Stationary vs. Difference Stationary

- **Trend Stationary:**

$$Y_t = \mu + \beta t + \epsilon_t$$

Detrend to achieve stationarity.
- **Difference Stationary:**

$$\Delta Y_t = Y_t - Y_{t-1}$$

Differencing removes non-stationarity.

ARIMA Models

- **ARIMA(p,d,q):**
 - p : Autoregressive lags.
 - d : Differencing order.
 - q : Moving average lags.

Support Vector Machines

- The objective of the support vector machine (SVM) algorithm is to find a hyperplane in a $(d - 1)$ -dimensional space (where d is the number of features) that distinctly classifies the data points.
- A "good"separator maximizes the margin between different classes.

Maximizing Margin

- Increasing the margin reduces model complexity, which helps in better generalization.
- Lesson from learning theory: If the hypothesis space H is constrained in size and/or the training dataset is large, low training error is likely to indicate low generalization error.

Linear Separators A linear separator is defined by a hyperplane: $\theta^T x = 0$. The decision function is given

by: $h(x) = \text{sign}(\theta^T x)$

SVM Optimization

- The primal problem is given by:

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2 \quad \text{s.t.} \quad y_i(\theta^T x_i) \geq 1, \quad \forall i$$

- The Lagrangian is used to transform the primal problem into a dual problem, which can be more efficient to solve.
- The dual formulation allows the use of kernel functions, which makes the problem computationally feasible in higher dimensions.

Kernel Trick

- SVMs can be extended to non-linear classification by using kernel functions to map the input features into a higher-dimensional space.
- Common kernels include:

– **Polynomial Kernel:** $K(x_i, x_j) = (\langle x_i, x_j \rangle + c)^d$

– **Gaussian Kernel (RBF):** $K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$

– **Sigmoid Kernel:** $K(x_i, x_j) = \tanh(\alpha x_i^T x_j + c)$

- The kernel trick allows SVMs to create complex decision boundaries without explicitly computing the higher-dimensional feature mappings.

Slack Variables and Soft Margin

- In cases where the data are not linearly separable, slack variables ξ_i are introduced to allow some misclassifications.
- The objective function is modified to minimize both the margin and the penalty for misclassified points:

$$\min_{\theta, \xi} \frac{1}{2} \sum_{j=1}^d \theta_j^2 + C \sum_{i=1}^n \xi_i \quad \text{s.t.} \quad y_i(\theta^T x_i) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

- The parameter C controls the trade-off between maximizing the margin and minimizing classification error.
- Strengths and Weaknesses of SVM**
- **Strengths:**

- Good generalization in both theory and practice.
- Works well with few training instances.
- Can find globally optimal solutions.
- Effective in high-dimensional spaces and when the number of dimensions exceeds the number of samples.

- **Weaknesses:**

- Slow to train/predict for large datasets.
- Requires careful selection of kernel and its parameters.
- Memory-intensive, especially with large datasets.
- Performance is sensitive to the choice of the regularization parameter C and the kernel parameters.

Applications of SVM

- SVMs are used in a wide range of applications, including:
 - extbfText classification: SVMs are effective in categorizing documents into topics.
 - extbfImage classification: SVMs can be used for recognizing and classifying images.
 - extbfBioinformatics: SVMs are used for protein classification and gene expression analysis.
 - extbfFinance: SVMs are applied in stock price prediction and risk management.

Tree Based Methods

Tree Based Methods

- Tree-based methods involve stratifying or segmenting the predictor space into a number of simple regions.
- These methods can be used for both regression and classification tasks.

Decision Trees

- Decision trees divide the feature space into axis-parallel (hyper-) rectangles.
- Each internal node represents a test on an attribute, each branch represents an outcome of the test, and each leaf node represents a class label or a continuous value (for regression).

Splitting Criteria

- Decision trees use splitting criteria such as Gini impurity, entropy, or mean squared error to determine the best attribute to split the data.
- extbfGini Impurity: Measures the impurity of a node; lower values indicate better splits.
- extbfEntropy: Measures the uncertainty or impurity; lower entropy after splitting indicates better separation of classes.

Bagging (Bootstrap Aggregation)

- Bagging is used to reduce the variance of decision trees by creating multiple bootstrapped datasets from the original training data and fitting a separate decision tree to each.
- The final prediction is made by averaging (for regression) or taking a majority vote (for classification).

Random Forests

- Random forests improve upon bagging by adding an element of randomness: at each split, a random subset of features is considered.
- This helps to reduce correlation between individual trees, leading to better model performance.

Boosting

- Boosting is an ensemble technique that builds trees sequentially, with each tree trying to correct the errors of the previous one.
- Common algorithms include AdaBoost and Gradient Boosting.

Strengths and Weaknesses of Tree-Based Methods

- **Strengths:**
 - Easy to interpret and visualize.
 - Can handle both numerical and categorical data.
 - Non-linear relationships between features do not affect tree performance.
- **Weaknesses:**
 - Prone to overfitting, especially with deep trees.
 - Unstable; small changes in the data can lead to different splits.

Applications of Tree-Based Methods

- Used in a variety of fields such as finance (credit scoring), healthcare (disease prediction), and marketing (customer segmentation).
- Random forests and boosting are often used in Kaggle competitions due to their high predictive power.

Neural Networks

Neural Networks

- Neural networks are a set of algorithms designed to recognize patterns, inspired by the structure and function of the human brain.
- They consist of layers of interconnected nodes (neurons) that can learn to approximate complex functions.

Perceptron

- The perceptron is the simplest type of artificial neural network, consisting of a single layer.
- It makes predictions based on a linear combination of input features and a step activation function.

Activation Functions

- Activation functions introduce non-linearity into the network, enabling it to learn complex patterns.
- Common activation functions include:
 - **Sigmoid:** $g(z) = \frac{1}{1+e^{-z}}$
 - **Tanh:** $g(z) = \tanh(z)$
 - **ReLU (Rectified Linear Unit):** $g(z) = \max(0, z)$
 - **Leaky ReLU:** $g(z) = \max(\alpha z, z)$, where α is a small constant

Feedforward Neural Network

- A feedforward neural network consists of an input layer, one or more hidden layers, and an output layer.
- Information flows in one direction, from input to output, without any feedback loops.

Backpropagation

- Backpropagation is the algorithm used to train neural networks by updating the weights to minimize the error.
- It involves calculating the gradient of the loss function with respect to each weight by using the chain rule.

Convolutional Neural Networks (CNNs)

- CNNs are specialized neural networks used for processing grid-like data, such as images.
- They consist of convolutional layers, pooling layers, and fully connected layers.
- Convolutional layers apply filters to detect features such as edges, textures, or patterns in the input.

Recurrent Neural Networks (RNNs)

- RNNs are designed for sequential data, where the output depends on previous elements in the sequence.
- They have loops that allow information to persist, making them suitable for tasks like time series prediction and natural language processing.
- Long Short-Term Memory (LSTM) networks are a type of RNN that can learn long-term dependencies.

Strengths and Weaknesses of Neural Networks

- **Strengths:**
 - Can learn complex non-linear relationships.
 - Highly flexible and adaptable to different types of data.
 - State-of-the-art performance on tasks such as image and speech recognition.
- **Weaknesses:**
 - Requires large amounts of data and computational resources.
 - Prone to overfitting, especially with deep architectures.
 - Difficult to interpret compared to simpler models.

Applications of Neural Networks

- Neural networks are used in a variety of applications, including:
 - **Image recognition:** Identifying objects within images.
 - **Natural language processing:** Language translation, sentiment analysis.
 - **Finance:** Stock price prediction, fraud detection.
 - **Healthcare:** Disease diagnosis, medical image analysis.