

Machine Learning Submodule

Model Assessment and Selection

Definition 1.1 Statistical Inference: Is the process of deducing properties of an underlying probability distribution by mere analysis of data.

Definition 1.2 Model Selection:

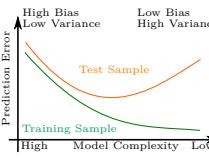
Is the process of selecting a model f from a given or chosen class of models \mathcal{F}

Definition 1.3 Hyperparameter Tuning: Is the process of choosing the hyperparameters θ of a given model $f \in \mathcal{F}$

Definition 1.4 Model Assessment/Evaluation: Is the process of evaluating the performance of a model.

Definition 1.5 Overfitting:

Describes the result of training/fitting a model f to closely to the training data $\mathcal{Z}^{\text{train}}$. That is, we are producing overly complicated model by fitting the model to the noise of the training set.



Consequences: the model will generalize poorly as the test set $\mathcal{Z}^{\text{test}}$ will not have not the same noise
⇒ big test error.

1. Empirical Risk Minimization

2. Generalization Error

Definition 1.6

Generalization/Prediction Error (Risk): Is defined as the expected value of a loss function l of a given predictor m , for data drawn from a distribution $P_{\mathcal{X}, \mathcal{Y}}$.

$$R_p(m) = \mathbb{E}_{(x,y) \sim p}[l(y; m(x))] = \int_D P(x, y) l(y; m(x)) dx dy \\ = \int_{\mathcal{X}} \int_{\mathcal{Y}} p(x, y) l(y, m(x)) dx dy \\ ?? \int_{\mathcal{X}} \int_{\mathcal{Y}} l(y, m(x)) p(y|x) p(x) dx dy \quad (1.1)$$

Interpretation

Is a measure of how accurately an algorithm is able to predict outcome values for future/unseen/test data.

Definition 1.7 Expected Conditional Risk: If we only know a certain x but not the distribution of those measurements ($x \sim P_{\mathcal{X}}(x)$), we can still calculate the expected risk given/conditioned on the known measurement x :

$$R_p(m, x) = \int_{\mathcal{Y}} l(y, m(x)) p(y|x) dy$$

Corollary 1.1 Note: $\stackrel{\text{def. 1.6}}{\iff} \stackrel{\text{def. 1.7}}{=}$:

$$R_p(m) = \mathbb{E}_{x \sim p}[R_p(m, x)] = \int_{\mathcal{X}} P(x) R_p(m, x) dx \quad (1.2)$$

1. Expected Risk Minimizer

Definition 1.8 Expected Risk Minimizer (TRM) m^* : Is the model m that minimizes the total expected risk:

$$m^* \in \arg \min_{m \in \mathcal{C}} R(m) = \arg \min_{m \in \mathcal{C}} \mathbb{E}_p[l(y; m(x))] \quad (1.3)$$

3. Empirical Risk

In practice we do neither know the distribution $P_{\mathcal{X}, \mathcal{Y}}(x, y)$, nor $P_{\mathcal{X}}(x)$ or $P_{\mathcal{Y}|X}(y|x)$ (otherwise we would already know the solution).

But: even though we do not know the distribution of $P_{\mathcal{X}, \mathcal{Y}}(x, y)$ we can still sample from it in order to define an empirical risk.

Definition 1.9 Empirical Risk:

Is the average of a loss function of an estimator h over a finite set of data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ drawn from $P_{\mathcal{X}, \mathcal{Y}}(x, y)$:

$$\hat{R}_n(m) = \frac{1}{n} \sum_{i=1}^n l(m(x_i), y_i)$$

1. Empirical Risk Minimizer

Definition 1.10 Empirical Risk Minimizer (ERM) \hat{m} : Is the model \hat{m} that minimizes the total empirical risk:

$$\hat{m} \in \arg \min_{m \in \mathcal{C}} \hat{R}(m) = \arg \min_{m \in \mathcal{C}} n^{-1} \sum_{i=1}^n l(m(x_i), y_i) \quad (1.4)$$

Questions

- ① How far is the true risk $R(m)$ from the empirical risk $\hat{R}(m)$, for a given m
- ② Given a chosen hypothesis class \mathcal{F} . How far is the minimizer of the true cost way from the minimizer of the empirical cost

$$m^*(x) \in \arg \min_{m \in \mathcal{F}} R(m) \quad \text{vs.} \quad \hat{m}(x) \in \arg \min_{m \in \mathcal{F}} \hat{R}(m)$$

We hope that $\lim_{n \rightarrow \infty} \hat{R}_n(m) = R(m)$.

3.1.1. Squared Loss Expected Squared Risk

Definition 1.11 Mean Squared Error (MSE):

$$R(m) = \text{MSE}(x) = \mathbb{E}[(\hat{m}(x) - m(x))^2] \quad (1.5)$$

Corollary 1.2 title:

$$\text{MSE}(x) = \text{Bias}^2(x) + \text{V}(x) = (\mathbb{E}[\hat{m}(x) - m(x)]^2 + \text{V}(\hat{m}(x))) \quad (1.6)$$

Definition 1.12

Integrated Means Squared Error (IMSE)/(MISE): the integrated MSE or *Mean integrated square error* (MISE) is defined as:

$$\text{IMSE} = \int_x \text{MSE}(x) dx = \int_x \mathbb{E}[(\hat{m}(x) - m(x))^2] dx \quad (1.7)$$

Empirical Squared Risk

Definition 1.13

Mean/Average Squared Prediction Error (MSPE): the empirical MSE or *Mean/Average Squared Error of Prediction* (MSEP)

$$\hat{R}_n(m) = \text{ave}_n(\hat{m})^2 = \frac{1}{n} \sum_{i=1}^n (\hat{m}(x_i) - m(x_i))^2 \quad (1.8)$$

Corollary 1.3

MSEP for new observations: Given a new observation x_{new} distributed as:

$$Y_{\text{new}} = m(x_{\text{new}}) + \epsilon \quad \epsilon \stackrel{\text{i.e.}}{\sim} \mathcal{N}(0, \sigma^2)$$

then it holds that:

$$\text{MSEP}(x_{\text{new}}) = \text{MSE}(x_{\text{new}}) + \sigma^2 \quad (1.9)$$

Explanation 1.1. The mean squared error of prediction does not go to zero if $n \rightarrow \infty$ as it has an irreducible noise σ .

Definition 1.14

[example 3.9], [proof 3.1]

Bayes' optimal predictor for the L2-Loss:

Assuming: i.i.d. generated data by $(x_i, y_i) \sim P_{\mathcal{X}, \mathcal{Y}}$.

Considering: the least squares risk:

$$R_p(h) = \mathbb{E}_{(x,y) \sim p}[(y - h(x))^2]$$

The best hypothesis/predictor h^* minimizing $R(h)$ is given by **conditional mean/expectation** of the data:

$$h^*(x) = \mathbb{E}[Y|X=x] \quad (1.10)$$

Cross Validation

Definition 1.15 Cross Validation: Is a model validation/assessment techniques in order to improve the model generalization performance.

Explanation 1.2. Cross validation helps to increase the model ability to predict out of sample data.

Definition 1.16 Labeled Data

$$\mathcal{D}/\mathcal{Z}: \quad \mathcal{Z} = \mathcal{D} := \{(z_j = (x_j, y_j) \mid x_j \in \mathcal{X}, y_j \in \mathcal{Y})\}$$

2. Training Set

Definition 1.17 Training Set

$\mathcal{Z}^{\text{train}} \subset \mathcal{Z}$: Is a part of the data on which we train our model \hat{m} in order to reduce the empirical

$$\mathcal{Z}^{\text{train}} = \{(x_1^{\text{train}}, y_1^{\text{train}}), \dots, (x_n^{\text{train}}, y_n^{\text{train}})\}$$

Definition 1.18

[??]
Training Error $\hat{R}(\hat{f}, \mathcal{Z}^{\text{train}})$: is the model that minimizes the empirical risk [def. 1.10] on the training data [def. 1.17]:

$$\hat{m} \in \arg \min_{\hat{m} \in \mathcal{F}} \hat{R}(\hat{m}, \mathcal{Z}^{\text{train}}) \quad (1.11)$$

$$= \arg \min_{\hat{m} \in \mathcal{F}} n^{-1} \sum_{(x_i, y_i) \in \mathcal{Z}^{\text{train}}} l(\hat{m}(x_i), y_i)$$

3. Testing Set

Definition 1.19

[??]
Test Set $\mathcal{Z}^{\text{test}} \subset \mathcal{Z}$: Is part of the data that is used in order to test the performance of our model.

$$\mathcal{Z}^{\text{test}} = \{(x_1^{\text{test}}, y_1^{\text{test}}), \dots, (x_m^{\text{test}}, y_m^{\text{test}})\}$$

Definition 1.20 Test Error

$\hat{R}(f, \mathcal{Z}^{\text{test}})$: Is the error over the test set $\mathcal{Z}^{\text{test}}$ of a predictor \hat{m} that has been trained on the training set [def. 1.17]:

$$\hat{R}(f, \mathcal{Z}^{\text{test}}) = n^{-1} \sum_{(x_i, y_i) \in \mathcal{Z}^{\text{test}}} l(\hat{m}(x_i), y_i) \quad (1.12)$$

4. Validation Set

Definition 1.21 Validation Set

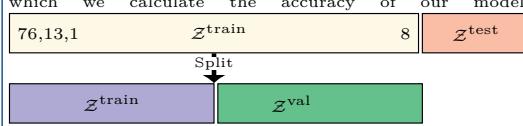
$\mathcal{Z}^{\text{val}} \subset \mathcal{Z}^{\text{train}}$: Is the part of the data that is used in order to select the our model \hat{m} from a given hypothesis class \mathcal{F} .

Explanation 1.3. We want to select a model \hat{m} from \mathcal{F} but in order to do so we need to determine the how well it predicts \Rightarrow validation set.

5. Validation Set/Split Once Approach

Definition 1.22 Hold out/Validation Set:

Split the data into a training set on which we train out model \hat{m} and a validation set on which we calculate the accuracy of our model:



Cons

- We do not use all information/data for training.
- We obtain a high variance estimate depending on the split.

Algorithm 1.1 Validation Set Approach:

Given: set of function classes \mathcal{F} and a loss l
1: train the model on the training set:

$$\hat{m} \in \arg \min_{m \in \mathcal{F}} \hat{R}(m, \mathcal{Z}^{\text{tr}}) = \arg \min_{m \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n l(y_i, m(x_i))$$

2: Determine the best parameter θ^* by using the validation set:
 $\hat{\theta}(\mathcal{Z}^{\text{val}}) \in \arg \min_{\theta: \hat{m}_{\theta} \in \mathcal{F}_{\theta}} \hat{R}(\hat{m}_{\theta}(\mathcal{Z}^{\text{tr}}), \mathcal{Z}^{\text{val}})$

3: Use the tests set in order to test the model:
 $\hat{R}(\hat{m}_{\theta(\mathcal{Z}^{\text{val}})}(\mathcal{Z}^{\text{tr}}), \mathcal{Z}^{\text{test}})$

Note: overfitting to the validation set

Tuning the configuration/hyperparameters of the model based on its performance on the validation set can result in overfitting to the validation set, even though your model is never directly trained on it \Rightarrow split the data into a test and training and validation set.

6. Leave-One-Out Cross Validation (LOOCV)

Definition 1.23

Leave One Out Cross-Validation (LOOCV):

Train n models on $n-1$ observations and use the left out observations for prediction:

$$\hat{m}_{n-1} \in \arg \min_{m \in \mathcal{F}} \frac{n-1}{n} \sum_{j=1, j \neq i}^n l(y_j, m(x_j)) \quad \forall i \in \{1, \dots, n\}$$

$$\hat{R}^{\text{LOOCV}} = n^{-1} \sum_{i=1}^n l(y_i, \hat{m}_{n-1}^{\text{LOOCV}}(x_i)) \quad (1.13)$$

Pros

- Is basically unbiased estimator, as we use $n-1$ training samples.
- Can have a high variance due to highly correlated training sets, as they only vary in one observation.
- Can be better as K -fold cross-validation for small data sets, as small data sets have usually a higher fluctuation \Rightarrow higher variance (as they are more sensitive to any noise/sampling artifacts).

Cons

- computational expensive, only for small data sets possible.
- Variance of the average can be very high due to highly correlated training sets.

3.6.1. LOOCV for Squared Loss and lin. Operator

Theorem 1.1 LOOCV Error for squared loss: For models that can be represented by a linear fitting operator S : $[\hat{m}(x_1), \dots, \hat{m}(x_n)]^T = SY$

$$n^{-1} \sum_{i=1}^n (y_i - \hat{m}_{n-1}^{\text{LOOCV}}(x_i))^2 = n^{-1} \sum_{i=1}^n \left(\frac{y_i - \hat{m}(x_i)}{1 - S_{ii}} \right)^2 \quad (1.14)$$

$$(1.15)$$

Definition 1.24 Generalized Cross Validation (GCV):

$$\text{GCV} = n^{-1} \sum_{i=1}^n \frac{(y_i - \hat{m}(x_i))^2}{(1 - n^{-1} \text{tr}(S))^2} \quad (1.16)$$

Explanation 1.4. It holds $\bar{S}_{ii} = \frac{1}{n} \sum_{i=1}^n S_{ii} = \frac{1}{n} \text{tr}(S)$ thus we can rewrite the mean as the trace, which can efficiently calculated in $\mathcal{O}(n)$.

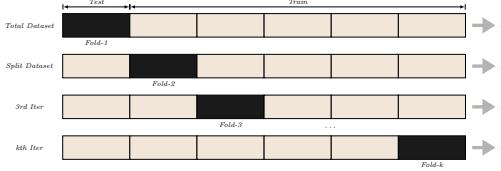
Note

GCV is a misnomer as it is an approximation and not a generalization.

7. K-Fold Cross Validation

Explanation 1.5 (*K*-fold Cross-Validation).

- 1 use all of the data by splitting the data into K random folds.
- 2 Calculate the training error K times by leaving out the k -th fold, fit the model to the other $K-1$ combined folds (training set) of size $n \cdot \frac{K-1}{K}$.
- 3 Do this by choosing each fold $k = 1, \dots, K$ once as validation set and calculate cross-validation error by averaging over them.



Definition 1.25 **K-fold Cross Validation:**

$$\mathcal{Z} = \mathcal{Z}_1 \cup \dots \cup \mathcal{Z}_\nu \cup \dots \cup \mathcal{Z}_K \quad \forall k \in \{1, \dots, K\}$$

$$\widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k} \in \arg \min_{m \in \mathcal{F}} \frac{|\mathcal{Z}_k|}{|\mathcal{Z}|} \sum_{i \in \mathcal{Z} \setminus \mathcal{Z}_k} l(y_i, m(x_i)) \quad (1.17)$$

$$\hat{\mathcal{R}}^{\text{CV}} = K^{-1} \sum_{k=1}^K |\mathcal{Z}_k|^{-1} \sum_{i \in \mathcal{Z}_k} l(y_i, \widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k}(x_i)) \quad (1.18)$$

Note

A good heuristic for choosing K is 5, or 10 or:
 $k = \min(\sqrt{n}, 10)$

Pros

- faster than LOOCV.

Cons

- runs $\approx K$ times slower than training/test-split, as we need to train the model K times.
- Has higher bias than LOOCV.
 There exists systematic tendency to underfit, as each of the K -fold cross validation models uses only $n \cdot \frac{K-1}{K}$ training samples
 ⇒ the estimates of prediction error will typically be more biased (towards simpler models), as the bias increases with a lower number of samples/d.o.f. (see Rao Cramer).
- Depends on the explicit realization of the K subsets.

8. Many Random Divisions

Definition 1.26 **Leave d -out CV:**

Generalize LOOCV/ d -fold CV by considering all possible realizations eq. (64.3) of d samples:

$$\mathcal{Z} = \mathcal{Z}_1 \cup \dots \cup \dots \cup \mathcal{Z}_{\binom{n}{d}} \quad \forall k \in \left\{1, \dots, \binom{n}{d}\right\}$$

$$\widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k} \in \arg \min_{m \in \mathcal{F}} \frac{|\mathcal{Z}_k|}{|\mathcal{Z}|} \sum_{i \in \mathcal{Z} \setminus \mathcal{Z}_k} l(y_i, m(x_i)) \quad (1.19)$$

$$\hat{\mathcal{R}}^{\text{CV}} = \binom{n}{d}^{-1} \sum_{k=1}^{\binom{n}{d}} |\mathcal{Z}_k|^{-1} \sum_{i \in \mathcal{Z}_k} l(y_i, \widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k}(x_i)) \quad (1.20)$$

Explanation 1.6. Is a generalization of LOOCV as it does not depend on the indexing in comparison to classical K -CV.

Pros

- has often a smaller variance.

Data Preprocessing

A Statistical Perspective

1. Information Theory

1. Information Content

Definition 3.1 Information (Claude Elwood Shannon):
Information is the resolution of uncertainty.

Amount of Information

The information gained by the realization of a coin tossed n-times should equal to the sum of the information of tossing a coin once n-times:

$$I(p_0 \cdot p_1 \cdots p_n) = I(p_0) + I(p_1) + \cdots + I(p_n)$$

⇒ can use the logarithm to satisfy this

Definition 3.2 Surprise/Self-Information/-Content:
Is a measure of the information of a realization x of a random variable $X \sim p$:

$$I_X(x) = \log\left(\frac{1}{p(X=x)}\right) = -\log p(X=x) \quad (3.1)$$

Explanation 3.1 (Definition 3.2).

$I(A)$ measures the number of possibilities for an event A to occur in bits:

$$I(A) = \log_2 (\# \text{possibilities for } A \text{ to happen})$$

Corollary 3.1 Units of the Shannon Entropy:
The Shannon entropy can be defined for different logarithms

	log	units
Base 2		Bits/Shannons
Natural		Nats
Base 10		Dits/Bans

Explanation 3.2. An uncertain event is much more informative than an expected/certain event:

$$\text{surprise/inf. content} = \begin{cases} \text{big} & \text{if } p_X(x) \text{ unlikely} \\ \text{small} & \text{if } p_X(x) \text{ likely} \end{cases}$$

2. Entropy

Information content deals with a single event. If we want to quantify the amount of uncertainty/information of a probability distribution, we need to take the expectation over the information content^[def. 3.2]:

Definition 3.3 Shannon Entropy [example 3.3]:
Is the expected amount of information of a random variable $X \sim p$:

$$H(p) = \mathbb{E}_X[I_X(x)] = \mathbb{E}_X\left[\log\left(\frac{1}{p_X(x)}\right)\right] = -\mathbb{E}_X[\log p_X(x)] \\ = -\sum_{i=1}^n p(x_i) \log p(x_i) \quad (3.2)$$

Definition 3.4 Differential/Continuous entropy:
Is the continuous version of the Shannon entropy^[def. 3.3]:

$$H(p) = \int_{x \sim p} -f(x) \log f(x) dx \quad (3.3)$$

Notes

- The Shannon entropy is maximized for uniform distributions
- People sometimes write $H(X)$ instead of $H(p)$ with the understanding that p is the distribution of X .

Property 3.1 Non negativity:

Entropy is always non-negative:
 $H(X) \geq 0$ if X is deterministic $H(X) = 0$ (3.4)

1.2.1. Conditional Entropy

Proposition 3.1 Conditioned Entropy $H(Y|X = x)$:
Let X and Y be two random variables with a conditional pdf $p_{Y|X}$. The entropy of Y conditioned on X taking a certain value x is given as:

$$H(Y|X = x) = \mathbb{E}_{Y|X=x}\left[\log \frac{1}{p_{Y|X}(Y|X=x)}\right] \\ = -\mathbb{E}_{Y|X=x}\left[\log p_{Y|X}(y|X=x)\right] \quad (3.5)$$

Definition 3.5 Conditional Entropy proof 3.4
 $H(Y|X)$:

Is the amount of information needed to determine Y if we already know X and is given by averaging $H(Y|X = x)$ over X :

$$H(Y|X) = [\mathbb{E}_X H(Y|X = x)] = -\mathbb{E}_{X,Y}\left[\log \frac{p(x,y)}{p(x)}\right] \quad (3.6) \\ = \mathbb{E}_{X,Y}\left[\log \frac{p(x)}{p(x,y)}\right]$$

Definition 3.6 Chain Rule for Entropy: proof 3.5

$$H(Y|X) = H(X, Y) - H(X) \\ H(X|Y) = H(X, Y) - H(Y) \quad (3.7)$$

Property 3.2 Monotonicity:

Information/conditioning reduces the entropy

$$\Rightarrow \text{Information never hurts.} \quad H(X|Y) \geq H(X) \quad (3.8)$$

Corollary 3.2 From eq. (3.17):
 $H(X, Y) \leq H(X) + H(Y)$ (3.9)

3. Cross Entropy

Definition 3.7 Cross Entropy [proof 3.3]:

Lets say a model follows a true distribution $X \sim p$ but we model X with a different distribution $X \sim q$. The cross entropy between p and q measure the average amount of information/bits needed to model an outcome $x \sim X \sim p$ with q :

$$H(p, q) = \mathbb{E}_{x \sim p}\left[\log\left(\frac{1}{q(x)}\right)\right] \quad (3.10)$$

$$= -\mathbb{E}_{x \sim p}[\log q(x)] \quad (3.11)$$

$$= H(p) + D_{KL}(p \parallel q) \quad (3.12)$$

Corollary 3.3 Kullback-Leibler Divergence:

$D_{KL}(p \parallel q)$ measures the extra price (bits) we need to pay for using q .

4. Kullback-Leibler (KL) divergence

If we want to measure how different two distributions q and p are w.r.t. to the same random variable X , we can define another measure.

Definition 3.8

Kullback-Leibler divergence. [examples 3.4 and 3.7]

/Relative Entropy from p to q : Given two probability distributions p, q of a random variable X . The Kullback-Leibler divergence is defined to be:

$$D_{KL}(p \parallel q) = \mathbb{E}_{x \sim p}\left[\log \frac{p(x)}{q(x)}\right] = \mathbb{E}_{x \sim p}[\log p(x) - \log q(x)] \quad (3.13)$$

and measures how far away a distribution q is from a another distribution p .

Explanation 3.3.

- p decides where we put the mass if $p(x)$ is zero we do not care about $q(x)$.
- $p(x)/q(x)$ determines how big the difference between the distributions is.

Intuition

The KL-divergence helps us to measure just how much information we lose when we choose an approximation.

1.2.2. Properties of KL-Divergence

Property 3.3 Non-Symmetric:

$$D_{KL}(p \parallel q) \neq D_{KL}(q \parallel p) \quad \forall p, q \quad (3.14)$$

Property 3.4:

$$D_{KL}(p \parallel q) \geq 0 \quad (3.15)$$

$$D_{KL}(p \parallel q) = 0 \iff p(x) = q(x) \forall x \in \mathcal{X} \quad (3.16)$$

Note

The KL-divergence is not a real distance measure as $KL(p \parallel Q) \neq KL(Q \parallel P)$

Corollary 3.4 Lower Bound on the Cross Entropy: The entropy provides a lower bound on the cross entropy, which follows directly eq. (3.16). from

5. Jensen-Shanon Divergence

6. Mutual Information

Definition 3.9

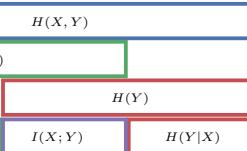
example 3.8

Mutual Information/Information Gain: Let X and Y be two random variables with a joint probability distribution. The mutual information of X and Y is the reduction in uncertainty in X if we know Y and vice versa.

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad (3.17)$$

$$= H(X) + H(Y) - H(X, Y)$$

$$= D_{KL}(p_{X,Y} \parallel p_X p_Y)$$



Explanation 3.4 (Definition 3.9).

$$I(X; Y) = \begin{cases} \text{big} & \text{if } X \text{ and } Y \text{ are highly dependent} \\ 0 & \text{if } X \text{ and } Y \text{ are independent} \end{cases} \quad (3.18)$$

Property 3.5 Symmetry:

$$I(X; Y) = I(Y; X)$$

Property 3.6 Positiveness:

$$I(X; Y) \geq 0 \quad \text{if } X \perp\!\!\!\perp Y \quad I(X; Y) = 0 \quad (3.19)$$

Property 3.7:

$$I(X; Y) \leq H(X) \quad I(X; Y) \leq H(Y) \quad (3.20)$$

Property 3.8 Self-Information:

$$H(X) = I(X; X)$$

Property 3.9 Montone Submodularity: Mutual information is monotone submodular^[def. 50.14]:

$$H(X, z) - H(z) \geq H(Y, z) - H(Y) \quad (3.21)$$

$$\Leftrightarrow H(z|X) \geq H(x|Y) \quad (3.22)$$

2. Proofs

Proof 3.1 Bayes Optimal Predictor^[def. 1.14]:

$$\begin{aligned} R(h) &= \min_h \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p} [(y - h(\mathbf{x}))^2] \\ &\stackrel{?}{=} \min_h \mathbb{E}_{\mathbf{x} \sim p_{\mathcal{X}}} \mathbb{E}_{\mathbf{y} \sim p_{\mathcal{Y}|\mathcal{X}}} [(y - h(\mathbf{x}))^2 | \mathbf{x}] \\ &\stackrel{?}{=} \mathbb{E}_{\mathbf{x} \sim p_{\mathcal{X}}} \left[\min_h \mathbb{E}_{\mathbf{y} \sim p_{\mathcal{Y}|\mathcal{X}}} [(y - h(\mathbf{x}))^2 | \mathbf{x}] \right] \end{aligned}$$

(def. 1.7)

Now lets minimize the conditional executed risk:

$$h^*(\mathbf{x}) = \arg \min_h \mathbb{E}_{\mathbf{y} \sim p_{\mathcal{Y}|\mathcal{X}}} [(y - h(\mathbf{x}))^2 | \mathbf{x}] \quad (3.23)$$

$$\begin{aligned} 0 &= \frac{d}{dh^*} \mathcal{R}_p(h^*, \mathbf{x}) = \frac{d}{dh^*} \int (y - h^*)^2 p(y|x) dy \\ &= \int \frac{d}{dh^*} (y - h^*)^2 p(y|x) dy = \int 2(y - h^*) p(y|x) dy \\ &= -2h^* \underbrace{\int p(y|x) dy}_{=1} + \underbrace{\int y p(y|x) dy}_{\mathbb{E}_Y[Y|X=x]} \end{aligned}$$

Proof 3.2 Irreducible Error^[cor. 1.3]:

$$\begin{aligned} \text{MSEP}(x_n) &= \mathbb{E}[(Y - \hat{Y}(x_n))^2] = \mathbb{E}[(Y - \hat{m}(x_n))^2] \\ &= \mathbb{E}[(\epsilon + m(x_n) - \hat{m}(x_n))^2] \\ &= \mathbb{E}[\epsilon^2] + 2\mathbb{E}[\epsilon \cdot (m(x_n) - \hat{m}(x_n))] \\ &\quad + \mathbb{E}[(m(x_n) - \hat{m}(x_n))^2] \\ &= \mathbb{E}[\epsilon^2] + 2\mathbb{E}[\epsilon \cdot (m(x_n) - \hat{m}(x_n))] \\ &\quad + \mathbb{E}[(m(x_n) - \hat{m}(x_n))^2] \\ &= \mathbb{V}[\epsilon] + 2\mathbb{E}[\epsilon] \cdot \mathbb{E}[(m(x_n) - \hat{m}(x_n))] \\ &\quad = 0 \\ &\quad + \mathbb{E}[(\epsilon + m(x_n) - \hat{m}(x_n))^2] \\ &= \mathbb{V}[\epsilon] + \text{MSE}(x_n) \end{aligned}$$

Proof 3.3: Cross Entropy^[def. 3.7]

$$\begin{aligned} \mathbb{E}_{x \sim q} \left[\log \left(\frac{1}{p(x)} \right) \right] &= \mathbb{E}_{x \sim q} \left[\log \left(\frac{1}{p(x)} \right) + \log \left(\frac{q(x)}{q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim q} \left[\log \left(\frac{q(x)}{p(x)} \right) + \log \left(\frac{1}{q(x)} \right) \right] \\ &= H(p) + D_{\text{KL}}(p \parallel q) \end{aligned}$$

Notes: ♡

Since we can pick $h(\mathbf{x}_i)$ independently from $h(\mathbf{x}_j)$.

Note

$$\begin{aligned} \mathbb{E}[X] \mathbb{E}[Y|X] &= \int_X p_X(x) dx \int_Y p(y|x) dy \\ &= \int_X \int_Y p_X(x) p(y|x) xy dx dy = \mathbb{E}[X, Y] \end{aligned}$$

Proof 3.4: Definition 3.5

$$\begin{aligned} \mathbb{E}_X [H(Y|X = x)] &= \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x) \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x)p(y|x) \log p(y|x) \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \left(\frac{p(x, y)}{p(x)} \right) \end{aligned}$$

Proof 3.5: ^[def. 3.6] We start from eq. (3.6):

$$\begin{aligned} H(Y|X) &= -\mathbb{E}_{X,Y} \left[\log \frac{p(x,y)}{p(x)} \right] \\ &= -\sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)} + \sum_x p(x) \log \frac{1}{p(X)} \\ &= H(X, Y) - H(X) \end{aligned}$$

Proof 3.6: example 3.4

$$\begin{aligned} \text{KL}(p \parallel q) &= \mathbb{E}_p [\log(p) - \log(q)] \\ &= \mathbb{E}_p \left[\frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2} (\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right. \\ &\quad \left. + \frac{1}{2} (\mathbf{x} - \mu_q)^T \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] \\ &= \frac{1}{2} \mathbb{E}_p \left[\log \frac{|\Sigma_q|}{|\Sigma_p|} \right] - \frac{1}{2} \mathbb{E}_p \left[(\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right] \\ &\quad + \frac{1}{2} \mathbb{E}_p \left[(\mathbf{x} - \mu_q)^T \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] \\ &= \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2} \mathbb{E}_p \left[(\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right] \\ &\quad + \frac{1}{2} \mathbb{E}_p \left[(\mathbf{x} - \mu_q)^T \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] \end{aligned}$$

$$\begin{aligned} \mathbb{E}_p[a] &= \mathbb{E}_p [\text{tr} \{ (\mathbf{x} - \mu_p)^T \Sigma_p^{-1} (\mathbf{x} - \mu_p) \}] \\ &\stackrel{\text{eq. (59.56)}}{=} \mathbb{E}_p [\text{tr} \{ (\mathbf{x} - \mu_p)(\mathbf{x} - \mu_p)^T \Sigma_p^{-1} \}] \\ &= \mathbb{E}_p [\text{tr} \{ \Sigma_p \Sigma_p^{-1} \}] \\ &\stackrel{\text{eq. (59.56)}}{=} \mathbb{E}_p [\text{tr} \{ I_d \}] = \mathbb{E}_p[d] = d \\ \mathbb{E}_p[b] &\stackrel{\text{eq. (65.54)}}{=} (\mu_p - \mu_q)^T \Sigma_q^{-1} (\mu_p - \mu_q) + \text{tr} \{ \Sigma_q^{-1} \Sigma_p \} \end{aligned}$$

3. Examples

Example 3.1 : Normal distribution has two population parameters: the mean μ and the variance σ^2 .

Example 3.2 Various kind of estimators:

- Best linear unbiased estimator (BLUE).
- Minimum-variance mean-unbiased estimator (MVUE): minimizes the risk (expected loss) of the squared-error loss-function.
- Minimum mean squared error (MMSE).
- Maximum likelihood estimator (MLE): is given by the least squares solution (minimum squared error), assuming that the noise is i.i.d. Gaussian with constant variance and will be considered in the next section.

Example 3.3 Entropy of a Gaussian:

$$\begin{aligned} H(\mathcal{N}(\mu, \Sigma)) &= \frac{1}{2} \ln |2\pi e \Sigma| \stackrel{\text{eq. (59.57)}}{=} \frac{1}{2} \ln ((2\pi e)^d |\Sigma|) \\ &= \frac{d}{2} \ln (2\pi e)^d + \log |\Sigma| \quad (3.24) \\ \Sigma &= \text{diag}(\sigma_1^2, \dots, \sigma_d^2) \quad \frac{1}{2} \ln |2\pi e| + \frac{1}{2} \sum_{i=1}^d \ln \sigma_i^2 \end{aligned}$$

Example 3.4 proof 3.6

KL Divergence of Gaussians:

Given two Gaussian distributions:

$$p = \mathcal{N}(\mu_p, \Sigma_p) \quad q = \mathcal{N}(\mu_q, \Sigma_q) \quad \text{it holds}$$

$$D_{\text{KL}}(p \parallel q) = \frac{\text{tr}(\Sigma_q^{-1} \Sigma_p) + (\mu_q - \mu_p)^T \Sigma_q^{-1} (\mu_q - \mu_p) - d + \ln \left(\frac{|\Sigma_q|}{|\Sigma_p|} \right)}{2}$$

Example 3.5 KL Divergence of Scalar Gaussians:

$$\theta \sim q(\theta | \lambda) = \mathcal{N}(\mu_q, \sigma_q^2) \quad \lambda = [\mu_q \ \sigma_q]$$

$$p = \mathcal{N}(\mu_p, \sigma_p^2)$$

$$D_{\text{KL}}(p \parallel q) = \frac{1}{2} \left(\frac{\sigma_p^2}{\sigma_q^2} (\mu_q - \mu_p)^2 \sigma_q^{-2} - 1 + \log \left(\frac{\sigma_q^2}{\sigma_p^2} \right) \right)$$

Example 3.6 KL Divergence of Diag. Gaussians:

$$\begin{aligned} \theta &\sim q(\theta | \lambda) = \mathcal{N}(\mu_q, \text{diag}(\sigma_1^2, \dots, \sigma_d^2)) \quad \lambda = [\mu_{1:d} \ \sigma_{1:d}] \\ p &= \mathcal{N}(\mu_p, \text{diag}(\sigma_1^2, \dots, \sigma_d^2)) \end{aligned}$$

Example 3.7 KL Divergence of Gaussians:

$$p = \mathcal{N}(\mu_p, \text{diag}(\sigma_1^2, \dots, \sigma_d^2)) \quad q = \mathcal{N}(0, I) \quad \text{it holds}$$

$$D_{\text{KL}}(p \parallel q) = \frac{1}{2} \sum_{i=1}^d (\sigma_i^2 + \mu_i^2 - 1 - \ln \sigma_i^2)$$

Example 3.8 Gaussian Mutual Information:

$$\begin{aligned} \text{Given } X &\sim \mathcal{N}(\mu, \Sigma) \quad Y = X + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I) \\ I(X; Y) &= H(Y) - H(Y|X) = H(Y) - H(\epsilon) \\ &\stackrel{\text{eq. (3.24)}}{=} \frac{1}{2} \ln (2\pi e)^d |\Sigma + \sigma^2 I| - \frac{1}{2} \ln (2\pi e)^d |\sigma^2 I| \\ &= \frac{1}{2} \ln \frac{(2\pi e)^d |\Sigma + \sigma^2 I|}{(2\pi e)^d |\sigma^2 I|} \\ &= \frac{1}{2} \ln |\Sigma + \sigma^2 I| \end{aligned}$$

Example 3.9 Bayes Optimal Predictor and MLE^[def. 1.14]: Problem: we do not know the real distribution $p_{\mathcal{Y}|\mathcal{X}}(y|\mathbf{x})$, which we need in order to find the bayes optimal predictor according to eq. (1.10).

Idea:

1. Use artificial data/density estimator $\hat{p}(\mathcal{Y}|\mathcal{X})$ in order to estimate $\mathbb{E}[\mathcal{Y}|\mathcal{X} = \mathbf{x}]$
2. Predict a test point \mathbf{x} by:

$$\hat{y} = \hat{\mathbb{E}}[\mathcal{Y}|\mathcal{X} = \mathbf{x}] = \int \hat{p}(y|\mathbf{X} = \mathbf{x}) y dy$$

Common approach: $p(\mathcal{X}, \mathcal{Y})$ may be some very complex (non-smooth, ...) distribution \Rightarrow need to make some assumptions in order to approximate $p(\mathcal{X}, \mathcal{Y})$ by $\hat{p}(\mathcal{X}, \mathcal{Y})$

Idea: choose parametric form $\hat{p}(Y|\mathbf{X}, \theta) = \hat{p}_\theta(Y|\mathbf{X})$ and then optimize the parameter θ which results in the so called maximum likelihood estimation section 1.

Supervised Learning

Definition 3.10 Statistical Inference: Goal of Inference

1. What is a good guess of the parameters of my model?
2. How do I quantify my uncertainty in the guess?

$$\mathcal{D} \xrightarrow{\text{Model Fitting}} (\mathcal{X} \xrightarrow{\text{Learning method}} \mathcal{Y}) \xrightarrow{\text{Prediction}} \hat{\mathbf{y}}$$

Recall: goal of supervised learning

Given: training data:

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \subseteq \mathcal{X} \times \mathcal{Y}$$

find a hypothesis $h: \mathcal{X} \mapsto \mathcal{Y}$ e.g.

$$h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$

$$h(\mathbf{x}) = \text{sing}(\mathbf{w}^\top \mathbf{x})$$

$$h(\mathbf{x}) = \sum_{i=1}^n \alpha_i \mathbf{k}(\mathbf{x}_i, \mathbf{x})$$

• Neural Networks (single hidden layer):

$$h(\mathbf{x}) = \sum_{i=1}^n \mathbf{w}_i^\top \phi(\mathbf{w}_i^\top \mathbf{x})$$

s.t. we minimize prediction error/empirical risk ^[def. 1.10].

Fundamental assumption

The data is generated i.i.d. from some unknown probability distribution:

$$(\mathbf{x}_i, y_i) \sim p_{\mathcal{X}, \mathcal{Y}}(\mathbf{x}_i, y_i)$$

Note

The distribution $p_{\mathcal{X}, \mathcal{Y}}$ is dedicated by nature and may be highly complex (not smooth, multimodal, ...).

4. Estimators

Definition 3.11 (Sample) Statistic: A statistic is a measurable function f that assigns a single value F to a sample of random variables:

$$\mathbf{X} = \{X_1, \dots, X_n\}$$

$$f: \mathbb{R}^n \mapsto \mathbb{R} \quad F = f(X_1, \dots, X_n)$$

E.g. F could be the mean, variance, ...

Note

The function itself is independent of the sample's distribution; that is, the function can be stated before realization of the data.

Definition 3.12 Statistical/Population Parameter: Is a parameter defining a family of probability distributions see example 3.1

Definition 3.13 (Point) Estimator $\hat{\theta} = \hat{\theta}(\mathbf{X})$:

Given: n-samples $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathcal{X}$ an estimator

$$\hat{\theta} = h(\mathbf{x}_1, \dots, \mathbf{x}_n) \quad (3.25)$$

is a statistic/random variable used to estimate a true (population) parameter θ ^[def. 3.12] see also example 3.2.

Note

The other kind of estimators are interval estimators which do not calculate a statistic **but** an interval of plausible values of an unknown population parameter θ .

The most prevalent forms of interval estimation are:

- Confidence intervals (frequentist method).
- Credible intervals (Bayesian method).

Generalized Linear Models (GLMs)

Definition 3.14 Generalized Linear Model (GLM):

$$\mu = \mathbb{E} [\mathbf{Y} | \mathbf{X}] = g^{-1} (\eta) \quad (3.26)$$

$$\eta = \sum_{j=0}^p \beta_{jm} X_j \quad (3.27)$$

$$g (\mathbb{E} [\mathbf{Y} | \mathbf{X}]) = \eta \quad (3.28)$$

Generalized Additive Models (GAMs)

Definition 3.15 Generalized Additive Models (GAMs):

$$sdf \quad (3.29)$$

Regression

Definition 4.1

Explanatory-/Indep.-/Predi.-/Variables/Covariates \mathbf{x} : Are the input variable(s) that we want to relate to the response variable(s)^[def. 4.2].

Definition 4.2

Response-/Dependent-/Variable(s) \mathbf{y} : Are the output quantities that we are interested in.

Definition 4.3 Coefficients β : Are the coefficients that we are seeking.

Definition 4.4 Regression: Is the process of finding a possible relationship via some coefficients β between response-variables \mathbf{x} and a predictor-variable(s) \mathbf{y} up to some error ϵ :

$$\mathbf{y} = f(\mathbf{x}, \beta) + \epsilon \quad (4.1)$$

Note

The term regression comes from the latin term "regressus" and means "to go back" to something. Historically the term was introduced by Galton, who discovered that given an outlier point, further observations will regress back to the mean. In particular he discovered that children of very tall/small people tend to be a smaller/larger.

Definition 4.5 Linear Regression: Refers to regression that is linear w.r.t. to the parameter vector β (but not necessarily the data):

$$\mathbf{y} = \beta^\top \phi(\mathbf{x}) + \epsilon \quad (4.2)$$

Linearity

Linearity is w.r.t. the coefficients β_j .

Thus a model with transformed non-linear predictor^[def. 4.1] variables is still called *linear*.

Definition 4.6 Residual

Let us consider n observations $\{\mathbf{x}_i, y_i\}_{i=1}^n$. The residual (error) is the deviation of the observed values from the predicted values:

$$r_i := e_i = \hat{y}_i - y_i = y_i - \hat{\beta}^\top \mathbf{x}_i \quad i = 1, \dots, n \quad (4.3)$$

Simple (linear) regression (SLR)

Definition 4.7 [example 4.1]
Simple Linear Regression: Is a linear regression^[def. 4.8] with only one explanatory variable^[def. 4.1]:

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i = 1, \dots, n \quad (4.4)$$

Multiple (linear) regression (MLR)

Definition 4.8 Multiple Linear Regression: Is a linear regression model with multiple $\{\beta_j\}_{j=1}^p$ explanatory^[def. 4.1] variables:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$$

$$= \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \epsilon_i = \beta^\top \mathbf{x}_i + \epsilon_i \quad i = 1, \dots, n \quad (4.13)$$

$$\begin{bmatrix} \mathbf{x} \\ \vdots \\ \mathbf{x} \end{bmatrix} \begin{bmatrix} \beta \\ \vdots \\ \beta \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ \vdots \\ \mathbf{y} \end{bmatrix} \quad \text{Design Matrix: } \mathbf{X} \in \mathbb{R}^{n \times (p+1)}$$

$$\mathbf{y} \in \mathbb{R}^n \quad \beta \in \mathbb{R}^{p+1} \quad (4.5)$$

Note

Eq. 4.8 is usually an over-determined system of linear equations i.e. we have more observations than predictor variables.

Multiple vs. Multivariate lin. Reg.

Multivariate linear regression is simply linear regression with multiple response variables and thus nothing else but a set of simple linear regression models that have the same types of explanatory variables.

Definition 4.9

[example 4.2] **Simple Linear Quadratic Regression:** Is a linear regression^[def. 4.8] with two explanatory variables^[def. 4.1] written as:

$$y_i = \beta_1 + \beta_2 x_i + \beta_3 x_i^2 + \epsilon_i \quad i = 1, \dots, n \quad (4.6)$$

0.0.1. Existence

Corollary 4.1 Existence:

$$\begin{aligned} & x_{11}\beta_1 + x_{12}\beta_2 + \dots + x_{1p}\beta_p & y_1 \\ & x_{21}\beta_1 + x_{22}\beta_2 + \dots + x_{2p}\beta_p & y_2 \\ & \vdots & \vdots \\ & x_{n1}\beta_1 + x_{n2}\beta_2 + \dots + x_{np}\beta_p & y_n \\ \iff & \mathbf{y} \in \mathfrak{N}(\mathbf{X}) & (4.8) \end{aligned}$$

1. Linear/Oldinary Least Squares (OLS)

Problem: for an over determined system $n > p$ (usually) $\mathbf{y} \in \mathfrak{N}(\mathbf{X})$ (in particular given round off errors) s.t. there exists no parameter vector β that solves^[def. 4.8].

Idea: try to find the next best solution by minimizing the residual(s)^[def. 4.6].

Definition 4.10 Residual Sum of Squares:

Is the sum of residuals^[def. 4.6]:

$$\text{RSS}(\beta) := \sum_{i=1}^n e_i^2 = \sum_{i=1}^n \|y_i - \hat{y}_i\|_2^2 \quad (4.9)$$

Definition 4.11 Least Squares Regression lsq(\mathbf{X}, \mathbf{y}):

Minimizes the residual sum of squares:

$$\begin{aligned} \hat{\beta} & \in \arg \min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 = \arg \min_{\beta} \|\mathbf{y} - \mathbf{u}\|_2^2 \\ & \quad \mathbf{u} \in \mathfrak{N}(\mathbf{X}) \\ & = \arg \min_{\beta} \|\mathbf{r}\|_2^2 = \sum_{i=1}^n \left(\sum_{j=1}^p x_{ij} \beta_j - y_i \right)^2 = \text{RSS}(\beta) \end{aligned} \quad (4.10)$$

Alternative Formulation

Sometimes people write eq. (4.10) as $\frac{1}{2} \arg \min_{\beta} \|\mathbf{r}\|_2^2$ which leads to the same solution eq. (54.63).

2. Maximum Likelihood Estimate

Ridge MLE

Proposition 4.1 (Gauss Markov Assumptions)

Assumptions for Linear Regression Model:

- The $\{\mathbf{x}_i\}_{i=1}^n$ are deterministic and measured without errors.
- The variance of the error terms is homoscedastic^[def. 69.22]:

$$\mathbb{V}[e_i] = \sigma^2 < \infty \quad \forall i \quad (4.11)$$
- The errors are uncorrelated:

$$\text{Cov}[e_i, e_j] = 0 \quad \forall i \neq j \quad (4.12)$$
- The errors are jointly normally distributed with mean 0 and constant variance σ^2 :

$$e_i \sim \mathcal{N}(0, \sigma^2) \quad \forall i = 1, \dots, n \iff e \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n) \quad (4.13)$$

Definition 4.12

Simple Linear Regression Log-Likelihood:

Assume: a linear model $\mathbf{y} = \mathbf{X}\beta + \epsilon$
with Gaussian noise $\epsilon \sim \mathcal{N}(0, \mathbf{I}\sigma^2)$

With: $\mu = \mathbb{E}_\epsilon[\mathbf{y}] = \mathbb{E}_\epsilon[\mathbf{X}\beta + \epsilon] = \mathbf{X}\beta + 0$

$\mathbb{V}_\epsilon[\mathbf{y}] = \mathbb{V}_\epsilon[\mathbf{X}\beta + \epsilon] = 0 + \mathbb{V}_\epsilon[\epsilon] = \mathbf{I}\sigma^2$

Thus: $\mathbf{Y}|\mathbf{X} \sim \mathcal{N}(\mathbf{X}\beta, \mathbf{I}\sigma^2)$

with: $\theta = (\beta^\top, \sigma^2)^\top \in \mathbb{R}^{p+1}$

$$l_n(\mathbf{y}|\mathbf{X}, \theta) \propto -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^\top \mathbf{x}_i)^2 = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2$$

$$\theta^* \in \arg \max_{\theta \in \mathbb{R}^{p+1}} l_n(\mathbf{y}|\mathbf{X}, \theta) = \arg \min_{\theta \in \mathbb{R}^{p+1}} -l_n(\mathbf{y}|\mathbf{X}, \theta) \quad (4.14)$$

1. The Normal Equation

Definition 4.13

[proof 4.4] **The Normal Equations:**

Is the equation we need to solve in order to solve eq. (4.10) or equivalently eq. (4.14) and is no longer an over determined system:

$$\begin{bmatrix} \mathbf{x}^\top \mathbf{x} \\ \vdots \\ \mathbf{x}^\top \mathbf{x} \end{bmatrix} \begin{bmatrix} \beta \\ \vdots \\ \beta \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ \vdots \\ \mathbf{y} \end{bmatrix} \quad \mathbf{x}^\top \mathbf{x} \in \mathbb{R}^{p \times p} \\ \beta \in \mathbb{R}^p \\ \mathbf{x}^\top \mathbf{x} \in \mathbb{R}^{p \times n} \\ \mathbf{y} \in \mathbb{R}^n \quad (4.15)$$

Geometric Interpretation

Corollary 4.2 Geometric Interpretation:

[proof 4.5]

We want to find $\arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$ which is equal to finding:

$$\arg \min_{\beta \in \mathbb{R}^p} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} \in \{\mathbf{y} : \beta \in \mathbb{R}^p\} = \mathfrak{N}(\mathbf{X})$$

but this minimum is equal to the orthogonal projection^[def. 59.22] of \mathbf{y} onto $\mathfrak{N}(\mathbf{X})$ i.e. the map:

$$\mathbf{y} \mapsto \hat{\mathbf{y}}$$

is the orthogonal projection of \mathbf{y} onto $\mathfrak{N}(\mathbf{X})$.

Corollary 4.3 Orthogonality of residuals

[proof 4.6]: Corollary 4.2 implies that the residuals are orthogonal w.r.t. to all the column vectors of \mathbf{X} :

$$\mathbf{r}^\top \mathbf{x}^{(j)} = 0 \quad \forall j = 1, \dots, p \quad (4.16)$$

2.1.1. The Least Squares Solution

Proposition 4.2 Least Squares Solution:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (4.17)$$

Note

\mathbf{X}^\dagger is the Moore-Penrose pseudo-inverse of the matrix \mathbf{X} .

2.1.2. Solving The Normal Equation

Cholesky Decomposition

Corollary 4.4 Computational Complexity: $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{w} \in \mathbb{R}^d$ with n , the number of observations and d , the number of equations/features/dimension of the problem.

Assume: $d \leq n$, that is we have an overdetermined system, more equations than unknowns.

- Compute regular matrix (Matrix Product):

$$\mathbf{C} := \mathbf{X}^\top \mathbf{X} \in \mathcal{O}(n \cdot d^2)$$
.

- Compute the r.h.s. vector (Matrix-Vector):

$$\mathbf{c} := \mathbf{X}^\top \mathbf{y} \in \mathbb{R}^d \in \mathcal{O}(nd)$$
.

- Solve s.p.d. LSE via Cholesky decomposition:

$$\mathbf{C}\mathbf{w} = \mathbf{c} \in \mathcal{O}(d^3)$$
.

Thus the total cost amounts to $\mathcal{O}(d^3 + nd^2)$.

Note: s.p.d. C and cholesky decomposition

Assume: \mathbf{X} has a trivial kernel $\Leftrightarrow \mathbf{X}^\top \mathbf{X}$ is invertible.

- Symmetric:** a transposed matrix times itself is symmetric $\Rightarrow \mathbf{C}$ is symmetric.
- Positive definite:**

$$\mathbf{w}^\top \mathbf{C}\mathbf{w} = \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} = \|\mathbf{X}\mathbf{w}\|^2 > 0 \quad \forall \mathbf{w} \neq 0$$

QR Decomposition

2.1.3. Simple Linear Regression Solution

[proof 4.4] **Linear Regression Solution:**

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad \text{with} \quad \Sigma^2 = (\mathbf{X}^\top \mathbf{X})^{-1} \quad (4.18)$$

$$\mathbf{P} = \mathbf{X}^\top \mathbf{y}$$

$$\Sigma^2 : \text{Variane-Covar. M.} \quad \mathbf{P} : \text{Inp./Outp. Covariance}$$

$$\text{Moore-Penrose pseudo-inverse: } \mathbf{X}^\dagger \quad \text{with} \quad \mathbf{X}^\dagger \mathbf{X} = \mathbf{I} \quad (4.19)$$

2.1.4. Making Predictions

Definition 4.15 P/H = $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top : \mathbf{y} \mapsto \hat{\mathbf{y}}$

Hat/Projection Matrix:
Is the matrix that projects the \mathbf{y} onto the $\hat{\mathbf{y}}$:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} =: \mathbf{Py} \quad (4.20)$$

Property 4.1 Symmetry: \mathbf{P} is trivially symmetric.

Property 4.2 Idem-potent $\mathbf{P}^2 = \mathbf{P}$: \mathbf{P} is idem-potent i.e. projecting multiple times by \mathbf{P} is the same as projecting once.

Property 4.3 Trace:

$$\text{tr}(\mathbf{P}) = \text{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \text{tr}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X})$$

$$= \text{tr}(\mathbf{I}_{p \times p}) = p$$

Corollary 4.5 $\mathbf{P} : \mathbb{R}^n \mapsto \mathcal{X} \subseteq \mathbb{R}^p$: From these three properties it follows that \mathbf{P} is an orthogonal projection onto a p -dim subspace.

Corollary 4.6 Residual Projection: The residual can be represented in terms of eq. (4.20):

$$\mathbf{r} = (\mathbf{I} - \mathbf{P})\mathbf{y} \quad (4.21)$$

it follows that $\mathbf{I} - \mathbf{P}$ is an orthogonal projection onto $(n - p)$ -dim subspace $\mathcal{X}^\perp = \mathbb{R}^n \setminus \mathcal{X}$.

Uniqueness

Theorem 4.1: Let $\mathbf{A} \in \mathbb{R}^{p,p}$, $p \geq p$ then it holds that:

$$\mathbb{N}(\mathbf{A}) = \mathbb{N}(\mathbf{A}^\top \mathbf{A}) \quad \mathbb{R}(\mathbf{A}^\top \mathbf{A}) = \mathbb{R}(\mathbf{A} \mathbf{A}^\top) \quad (4.22)$$

Theorem 4.2 Full-Rank Condition **F.R.C.:**
Equation 4.13 has a unique least squares solution given by:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (4.23)$$

$$\Leftrightarrow \mathbb{N}(\mathbf{X}) = \{0\} \Leftrightarrow \text{rank}(\mathbf{X}) = p \quad p \geq p \quad (4.24)$$

2. Moments and Distributions

Property 4.4 Moments of $\hat{\beta}$ [proof 4.7]:

$$\mathbb{E}[\hat{\beta}] = \beta \quad \mathbb{V}[\hat{\beta}] = \text{Cov}[\hat{\beta}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \quad (4.25)$$

$$\hat{\beta} \sim \mathcal{N}_p(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}) \quad (4.26)$$

Property 4.5 Moments of $\hat{\mathbf{y}}$ [proof 4.9]:

$$\mathbb{E}[\hat{\mathbf{y}}] = \mathbb{E}[\mathbf{y}] = \mathbf{X}\beta \quad \mathbb{V}[\hat{\mathbf{y}}] = \text{Cov}[\hat{\mathbf{y}}] = \sigma^2 \mathbf{P} \quad (4.27)$$

$$\hat{\mathbf{y}} \sim \mathcal{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{P}) \quad (4.28)$$

Property 4.6 Moments of \mathbf{r} :

$$\mathbb{E}[\mathbf{r}] = 0 \quad \text{Cov}[\mathbf{r}] = \sigma^2 (\mathbf{I} - \mathbf{P}) \quad (4.29)$$

$$\mathbf{r} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{P})) \quad (4.30)$$

Property 4.7 Moments of $\hat{\sigma}^2$:

$$\hat{\sigma}^2 := \frac{1}{n-p} \sum_{i=1}^n r_i^2 \quad \Rightarrow \quad \mathbb{E}[\hat{\sigma}^2] = \sigma^2 \quad (4.31)$$

$$\hat{\sigma}^2 \sim \frac{\sigma}{n-p} \chi_{n-p}^2 \quad (4.32)$$

Note

The standard deviation σ^2 is given by $\epsilon \sim \mathcal{N}(0, \sigma^2)$. However we may not know σ^2 , thus we can estimate it by using the residuals \mathbf{r} .

Proof 4.1 Property 4.7: $\hat{\sigma}^2$ is an unbiased estimator of σ^2 :

2.2.1. The Gauss Markov Theorem

Theorem 4.3 Gauss–Markov theorem [proof 4.10]:

The BLUE of the β coefficients, of a linear regression model, satisfying the **Gauss–Markov assumptions** is given by the ordinary least squares (OLS) estimator, provided it exists (is invertible).

$$\mathbb{V}[\hat{\beta}] \leq \mathbb{V}[\tilde{\beta}] \quad \text{with} \quad \hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{C}\mathbf{y} \quad (4.33)$$

$\tilde{\beta}$ any lin. unb. est. for β

3. MLE with linear Model & Gaussian Noise

1. MLE for conditional linear Gaussians

Questions: what is $P(Y|X)$ if we assume a relationship of the form: We can use the MLE to estimate the parameters $\theta \in \mathbb{R}^k$ of a model/distribution h s.t.

$$y \approx h(\mathbf{x}; \theta) \iff y = h(\mathbf{x}; \theta) + \epsilon$$

\mathbf{x} : set of explicative variables. ϵ : noise/error term.

Lemma 4.1 : The conditional distribution D of Y given \mathbf{X} is equivalent to the unconditional distribution of the noise ϵ : $P(Y|\mathbf{X}) \sim D \iff \epsilon \sim D$

Example: Conditional linear Gaussian

Assume: a linear model $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ and Gaussian noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$

With $E[\epsilon] = 0$ and $y_i = \mathbf{w}^\top \mathbf{x}_i + \epsilon$, as well as ?? it follows: $y \sim \hat{P}(Y|\mathbf{X} = \mathbf{x}, \theta) \sim \mathcal{N}(\mu = h(\mathbf{x}), \sigma^2)$

with: $\theta = (\mathbf{w}^\top \sigma)^\top \in \mathbb{R}^{n+1}$

Hence Y is distributed as a linear transformation of the \mathbf{X} variable plus some Gaussian noise ϵ : $y_i \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \sigma^2) \Rightarrow$ Conditional linear Gaussian.

if we consider an i.i.d. sample $\{y_i, \mathbf{x}_i\}_{i=1}^n$, the corresponding conditional (log-)likelihood is defined to be:

$$\begin{aligned} L_n(Y|\mathbf{X}, \theta) &= \hat{P}(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \theta) \\ &\stackrel{\text{i.i.d.}}{=} \prod_{i=1}^n \hat{P}(Y|\mathbf{x}_i | y_i, \theta) = \prod_{i=1}^n \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} \exp\left(-\frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}\right) \\ &= (\sigma^2 2\pi)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2\right) \end{aligned}$$

$$L_n(Y|\mathbf{X}, \theta) = -\frac{n}{2} \ln \sigma^2 - \frac{n}{2} \ln 2\pi - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

$$\theta^* = \arg \max_{\mathbf{w} \in \mathbb{R}^d, \sigma^2 \in \mathbb{R}_+} L_n(Y|\mathbf{X}, \theta)$$

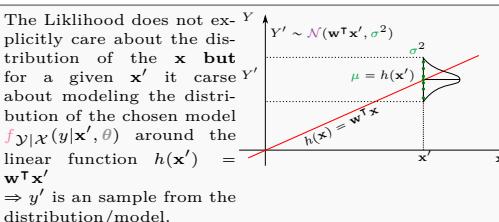
$$\frac{\partial \ln(Y|\mathbf{X}, \theta)}{\partial \theta} = \begin{pmatrix} \frac{\partial \ln(Y|\mathbf{X}, \theta)}{\partial w_1} \\ \vdots \\ \frac{\partial \ln(Y|\mathbf{X}, \theta)}{\partial w_d} \\ \frac{\partial \ln(Y|\mathbf{X}, \theta)}{\partial \sigma^2} \end{pmatrix} = \begin{pmatrix} \mathbf{0}_d \\ \vdots \\ \mathbf{0}_d \end{pmatrix}$$

$$\begin{aligned} \frac{\partial \ln(Y|\mathbf{X}, \theta)}{\partial \mathbf{w}} &= \frac{1}{\sigma^2} \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{w}^\top \mathbf{x}_i) = \mathbf{0} \in \mathbb{R}^d \\ &= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{w} = \sum_{i=1}^n \mathbf{x}_i y_i \\ \frac{\partial \ln(Y|\mathbf{X}, \theta)}{\partial \sigma^2} &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i - \mu)^2 = 0 \end{aligned}$$

$$\theta^* = \begin{pmatrix} \mathbf{w}^* \\ \sigma^* \end{pmatrix} = \begin{pmatrix} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i y_i \right) \\ \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^* \mathbf{x}_i)^2 \end{pmatrix} \quad (4.34)$$

Note

- The mean μ of the normal distribution follows from: $E[\mathbf{w}^\top \mathbf{x}_i + \epsilon_i] = E[\mathbf{w}^\top \mathbf{x}_i] + E[\epsilon_i] = \mathbf{w}^\top \mathbf{x}_i$ const. $= 0$
- The noise ϵ must have zero mean, otherwise it wouldn't be random anymore.
- The optimal function $h^*(\mathbf{x})$ determines the mean μ .
- We can also minimize: $\theta^* = \arg \max_{\theta} \hat{P}(Y|\mathbf{X}, \theta) = \arg \min_{\theta} -\hat{P}(Y|\mathbf{X}, \theta)$



Ridge Max Prior

Prior

Assume: prior $P(\beta|\Sigma)$ on the model parameter β is gaussian as well and depends on the hyperparameter $(\text{def. 6.7}) \Sigma \cong \text{co-variance matrix}$:

$$\begin{aligned} \beta \sim P_{\text{Ridge}}(\beta|\Sigma) &= \mathcal{N}(\beta|0, \Sigma) \\ &= (2\pi)^{-\frac{d+1}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \beta^\top \Sigma^{-1} \beta\right) \end{aligned}$$

$$\ln(\beta|\Sigma) = -\frac{1}{2} \ln \det(\Sigma)^{-1} - \frac{d+1}{2} \ln 2\pi - \frac{1}{2} \beta^\top \Sigma^{-1} \beta \quad (4.37)$$

Max Prior

$$\begin{aligned} \beta^* &\in \arg \max_{\beta \in \mathbb{R}^{d+1}} \ln(\beta|\Sigma) \\ &= \arg \max_{\beta \in \mathbb{R}^{d+1}} -\frac{1}{2} \ln \det(\Sigma)^{-1} - \frac{d+1}{2} \ln 2\pi - \frac{1}{2} \beta^\top \Sigma^{-1} \beta \end{aligned}$$

$$0 \stackrel{!}{=} \frac{\partial}{\partial \beta^*} \ln(\beta^*|\Sigma) = -\frac{\partial}{\partial \beta^*} \beta^* \Sigma^{-1} \beta^* \stackrel{\text{eq. (4.46)}}{=} -2 \Sigma^{-1} \beta^*$$

$$\begin{aligned} \beta^* &\in \arg \max_{\beta \in \mathbb{R}^{d+1}} \log(p(\beta|\Sigma)) = \arg \min_{\beta \in \mathbb{R}^{d+1}} -\ln(\beta|\Sigma) = 2 \Sigma^{-1} \beta^* \\ &\quad \beta \in \mathbb{R}^{d+1} \end{aligned} \quad (4.38)$$

Log-MAP

$$\begin{aligned} \beta^* &\in \arg \max_{\beta \in \mathbb{R}^{d+1}} P(\beta|\mathbf{X}, \mathbf{y}) \\ &= \arg \min_{\beta \in \mathbb{R}^{d+1}} -\log(P(\beta|\Sigma)) - \log(P(\mathbf{X}, \mathbf{y}|\beta)) \\ &\stackrel{\text{eq. (4.38)}}{=} \Sigma^{-1} \beta^* - \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y} + \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} \beta^* = 0 \\ &\iff (\Sigma^{-1} + \mathbf{X}^\top \mathbf{X} \sigma^{-2}) \beta^* = \sigma^{-2} \mathbf{X}^\top \mathbf{y} \\ &\quad (\sigma^2 \Sigma^{-1} + \mathbf{X}^\top \mathbf{X}) \hat{\beta} = \mathbf{X}^\top \mathbf{y} \\ &\hat{\beta}^{\text{MAP}} = (\sigma^2 \Sigma^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

Definition 4.16 Ridge MAP: For ridge regression we assume that the noise of the prior is uncorrelated/diagonal i.e.

$$\Sigma^{-1} = \mathbf{I} \sigma^{-2} \quad \text{and let} \quad \Lambda := \sigma^2 \Sigma^{-1} = \mathbf{I} \frac{\sigma^2}{\sigma^2} \quad (4.39)$$

which leads to:

$$\hat{\beta}^{\text{MAP}} = (\Lambda + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad \text{with} \quad \Lambda = \mathbf{I} \lambda = \mathbf{I} \frac{\sigma^2}{\sigma^2} \quad (4.40)$$

Definition 4.17 Regularization: Regularization is the process of introducing additional information/bias in order to solve an ill-posed problem or to prevent overfitting. (It is not feature selection)

Definition 4.18 Tikhonov regularization: Commonly used method of regularization of ill-posed problems.

$$\|\mathbf{X}\beta - \mathbf{y}\|^2 + \|\Gamma\beta\|^2 \quad (4.41)$$

G: Tikhonov matrix in many cases, this matrix is chosen as $\Gamma = \alpha \mathbf{I}$ giving preference to solutions with smaller norms; this is known as **Ridge/L2 regularization**.

Note

$$\frac{\partial \mathbf{w}^\top \mathbf{x}}{\partial \mathbf{w}} = \frac{\partial \mathbf{x}^\top \mathbf{w}}{\partial \mathbf{w}} = \mathbf{x}$$

3. MLE for general conditional Gaussians

Suppose we do not just want to fit linear functions but a general class of models $H_{sp} := \{h : \mathcal{X} \mapsto \mathbb{R}\}$ e.g. neural networks, kernel functions,...

Given: data $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ The MLE for general models h and i.i.d. Gaussian noise:

$$h \sim \hat{P}(Y|\mathbf{X}) = \hat{P}(Y|\mathbf{x} = \mathbf{x}, \theta) = \mathcal{N}(y|h^*(\mathbf{x}), \sigma^2)$$

Is given by the least squares solution:

$$h^* = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^n (y_i - h(\mathbf{x}_i))^2$$

E.g. for linear models $\mathcal{H} = \{h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \text{ with parameter } \mathbf{w}\}$

Other distributions

If we use other distributions instead of Gaussian noise, we obtain other loss functions e.g. L1-Norm for **Poisson Distribution**.

\Rightarrow if we know something about the distribution of the data we know which loss function we should choose.

Gaussian Prior/Likelihood MAP inference

$$\begin{aligned} \hat{\beta}^{\text{Ridge}} &= \arg \min_{\beta} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^\top \Lambda \beta \right\} \\ &= \arg \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \beta^\top \Lambda \beta \right\} \\ &\stackrel{\text{eq. (4.39)}}{=} \arg \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2 \right\} \\ &= \arg \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{i=1}^d \beta_i^2 \right\} \end{aligned}$$

$\|\mathbf{y} - \mathbf{X}\beta\|^2$ is forced to be small so that we find a weight vector β that matches the data as close as possible:

$$y_i = \beta_i \mathbf{x}_i + \epsilon_i \quad \text{s.t.} \quad \sum_{i=1}^n \epsilon_i \text{ small}$$

In other words we want to fit the data well.

- $\beta^\top \Lambda \beta \stackrel{\text{ridge}}{=} \lambda \|\beta\|^2$ says chose a model with a small magnitude $\|\beta\|^2$. Thus the smaller λ the bigger can the data faithfullness term be $\|\mathbf{y} - \mathbf{X}\beta\|^2$.

Note

The intercept β_0 in the regularizer term has to be left out. Penalization of the intercept would make the procedure depend on the origin chosen for y .

Thus we actually have (for data with non-zero mean):

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^d} \left\{ \|\mathbf{y} - (\mathbf{X}\beta + \beta_0)\|^2 + \lambda \sum_{i=1}^d \beta_i^2 \right\}$$

Note: SVD

Using SVD one can show that ridge regression shrinks first the eigenvectors with minimum explanatory variance. Hence L2/Ridge regression can be used to estimate the predictor importance and penalize predictors that are not important (have small explanatory variance).

Note: no feature selection

The coefficients in a ridge will go to zero as λ increases but will no become zero (as long as $\lambda \neq \infty$)! They are fit in a restricted fashion controlled by the **shrinkage penalty** λ .

$$\text{dofs}(\lambda) = \begin{cases} d & \text{if } \lambda = 0 \text{ (no regularization)} \\ \rightarrow 0 & \text{if } \lambda \rightarrow \infty \end{cases} \quad (4.42)$$

\Rightarrow Ridge cannot be used for variable selection since it retains all the predictors

Balance of $\lambda = \frac{\sigma^2}{\sigma^2}$ controls the tradeoff between simplicity and data faithfullness because:

- $\lambda \stackrel{\sigma \uparrow}{\rightarrow} \infty: \|\beta\|^2$ must be minimized:
 - $\sigma \uparrow$: model does not need to match data so perfectly as we have more noise in our data/observations \Rightarrow bigger errors (recall $\epsilon \sim \mathcal{N}(0, \mathbf{I}\sigma^2)$).
 - $\sigma \downarrow$: prior has smaller variance, thus our prior knowledge of the model is pretty exact/important (recall $\beta \sim \mathcal{N}(\beta|0, \mathbf{I}\sigma^2)$)
- $\lambda \stackrel{\sigma \downarrow}{\rightarrow} 0: \|\mathbf{y} - \mathbf{X}\beta\|^2$ must be minimized: model must match data perfectly
 - $\sigma \downarrow$: model does need to match perfectly, our observation/data has small variance/is well defined \Rightarrow do not allow big errors (recall $\epsilon \sim \mathcal{N}(0, \mathbf{I}\sigma^2)$).
 - $\sigma \uparrow$: our knowledge about the model is pretty vague (recall $\beta \sim \mathcal{N}(\beta|0, \mathbf{I}\sigma)$)

Note

- Often $\Lambda^{-1} = \mathbf{I} \in \mathbb{R}^{d+1 \times d+1}$

- Λ is symmetric and diagonal.

- $(d+1)$ dimension as we included offset into β .

Heuristic Map Inference

A really large weight vector β will result in amplifying noise/larger variance/fluctuations $\hat{=}$ overfitting.
This is because the complexity of the estimate increases with the magnitude of the parameter as it becomes easier to fit complex noise.

Ill-posed problem/Invertability and Ridge

Another advantage of Ridge regression is that, even if $\mathbf{X}^\top \mathbf{X}$ in eq. (4.40) is not invertible/regular/has not full rank. Then $(\mathbf{X}^\top \mathbf{X} + \Delta)$ will still be invertible/well posed. This was the original reason for L2/Ridge Regression.

$$\text{MAP} \hat{=} \text{Ridge}$$

$$\arg \max_{\mathbf{w}} P(\mathbf{w}|\mathbf{x}, y) = \arg \min_{\mathbf{w}} \lambda \|\mathbf{w}\|^2 + \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

MAP with a linear model and Gaussian noise equals classical ridge regression ??.

$$\begin{aligned} \arg \min_{\mathbf{w}} \lambda \|\mathbf{w}\|^2 + \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 &\equiv \arg \max_{\mathbf{w}} P(\mathbf{w}) \prod_{i=1}^n P(y_i|\mathbf{x}_i, \mathbf{w}) \\ \text{Ridge Regression} & \quad \text{MAP} \end{aligned}$$

Thus if we know our data β, σ we can chose λ statistically and do not need cross-validation.

Generalization

Regularized estimation can often be understood as MAP inference:

$$\begin{aligned} \arg \min_{\mathbf{w}} \sum_{i=1}^n l(\mathbf{w}^\top \mathbf{x}_i; \mathbf{x}_i, y_i) + C(\mathbf{w}) &= \\ &= \arg \max_{\mathbf{w}} \prod_{i=1}^n P(\mathbf{w}) P(y_i|\mathbf{x}_i, \mathbf{w}) = \arg \max_{\mathbf{w}} P(\mathbf{w}|\text{data}) \end{aligned}$$

$$\text{with } C(\mathbf{w}) = -\log P(\mathbf{w}) \quad l(\mathbf{w}^\top \mathbf{x}_i; \mathbf{x}_i, y_i) = -\log P(y_i|\mathbf{x}_i, \mathbf{w})$$

Priors

4. Laplace Prior $\hat{=}$ Lasso/L1-regularization

Intro

Question: what if $d \gg n$ e.g.

- bag of words with $d = \text{nb. of words} \gg n = \text{nb. of documents}$.
- Genome analysis $d = \text{nb. of genes} \gg n = \text{patients}$.

Problem: we have more unknowns/parameters than observations \Rightarrow no unique solution. **e.g.:** Trying to fit 1 data point with polynomial of degree 12.

Question: can we somehow still find a good solution if $n = \mathcal{O}(\ln d) \iff \exp. \text{more dim. than observations}$

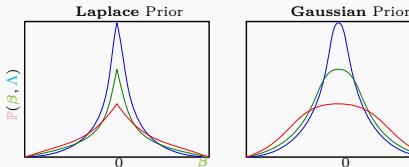
Idea: If most of the dimensions are irrelevant for the problem, then we can find a good (**sparse**) solution $\hat{=}$ **feature selection/dimensionality reduction**.

Given: Laplacian model prior $\beta \sim p(\beta|\Lambda)$:

$$\text{P-Lasso } (\beta|\Lambda) \text{ eq. (66.58)} \frac{\Lambda}{2} e^{-\Lambda |\beta|} = \prod_{j=1}^d \frac{\lambda_j}{2} e^{-\lambda_j |\beta_j|}$$

With $\Lambda^{-1} := \Sigma$ hyperparameter/covariance matrix

This leads to a L1 regularized model:



Thus: laplace priors gives sparseness, higher liklihood to get value at $\beta = 0$.

$$-\ln P(\beta|\Lambda) = \sum_{j=1}^d \lambda_j |\beta_j| - d \ln \frac{\lambda_j}{2} \quad (4.43)$$

Laplacian MAP Prior Inference

$$\begin{aligned} \beta^* &= \arg \min_{\beta \in \mathbb{R}^d} \left\{ \|y - (\mathbf{X}\beta + \beta_0)\|^2 + \lambda \|\beta\|_1 \right\} \\ &= \arg \min_{\beta \in \mathbb{R}^d} \left\{ \|y - (\mathbf{X}\beta + \beta_0)\|^2 + \lambda \sum_{i=1}^d |\beta_i| \right\} \quad (4.44) \end{aligned}$$

$|\beta_i|$ does not change β_i while β_i^2 becomes very small for values $\in (0, 1)$ thus when minimizing the L2 error $\|\beta\|_2^2 \rightarrow 0$ but not β_i while for L1 regularization will actually have to set β_i values to zero for large enough λ .

Advantage

Combines advantages of Ridge regression (convex function/optimization) and L0-regression (sparse and easy to interpret solution).

Difference L1& L2 penalties

Typically ridge or L2 penalties are much better for minimizing prediction error rather than L1 penalties. The reason for this is that when two predictors are highly correlated, L1 regularizer will simply pick one of the two predictors. In contrast, the L2 regularizer will keep both of them and jointly shrink the corresponding coefficients a little bit. Thus, while the L1 penalty can certainly reduce overfitting, you may also experience a loss in predictive power.

Notes

The unconstrained convex (see [cor. 54.12]) optimization problem eq. (4.44) is not differentiable at $\beta_i = 0$ and thus has no closed form solution as the L2 problem \Rightarrow quadratic programming.

5. Sparseness Priors/L0-regularization

$$-\ln P(\beta|s) = s \sum_{j=1}^d \mathbb{1}_{\beta_j \neq 0} = s \sum_{j=1}^d \begin{cases} 1 & \text{if } \beta_j \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (4.45)$$

\Rightarrow measure for the number of possible non-zero dimensions/parameters in β .

Advantage

- Leads always to sparse solution.
- Indicates/Explains model well as we only get a few non-zero parameters that determine/characterize the model.

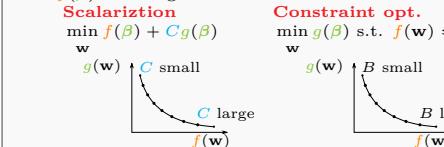
Drawback

Non-convex, non-differentiable problem \Rightarrow computationally difficult combinatorics.

Scalarization vs. Constrained Optimization

There are two equivalent ways of trading:

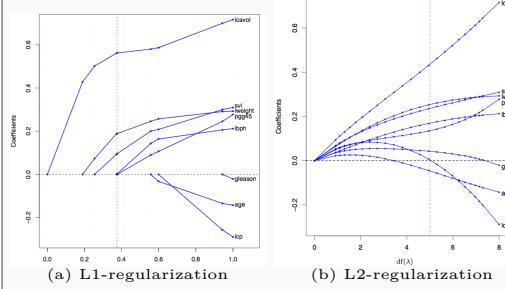
- $g(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2$: the data term and
- $f(\beta)$: the Regularizer.



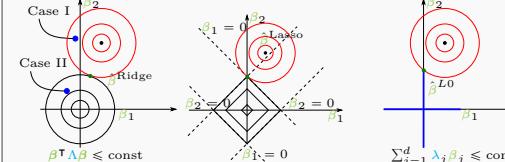
Note

Scalarization and **constrained optimization** gives the same curves $\iff f, g$ are both convex functions.
This is not necessarily for the same values of C and B but their exists always a relationship $C = u(B)$ s.t. this is true.

Comparison of priors



The constraint formulation of the optimization problems can be plotted for two features β_1, β_2 as:



- **Ridge Regression/L2-regression:** if the least squares error solution satisfies the constraint, we are fine (Case II), otherwise we do violated the constraint $\beta_1^2 + \beta_2^2 \leq \text{const}$ (Case I).
- **Lasso/L1-regression:** Here the constraint equals $|\beta_1| + |\beta_2| \leq \text{const}$ and leads to polyhedron. Most of the time we obtain a sparse solution $\hat{=}$ corner, due to the fact that corner regions increases much faster in volume, as the mixed regions (sparseness increases with number of dimensions).
- **Sparseness prior/L0-regression:** Leads to a super spiky geometry \Rightarrow always leads to a sparse solution.

Likelihoods

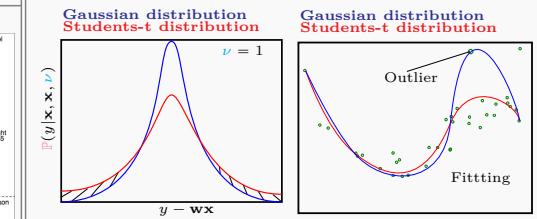
6. Student's-t likelihood loss function

$$\text{Students-t Distribution: } f(y|\mathbf{x}, \mathbf{w}, \nu, \sigma^2) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi \nu \sigma^2} \Gamma(\frac{\nu}{2})} \left(1 + \frac{(y - \mathbf{w}^\top \mathbf{x})^2}{\nu \sigma^2} \right)^{-\frac{\nu+1}{2}}$$

ν : determines speed of decay.

Problem L2/squared loss functions lead to estimates that are sensitive to outliers, that is because something that is far away from the expected value, will be increased/influences the model very much.

- For **Gaussian noise**: outliers are very unlikely and thus will have a big influence on the model.
- For **Students-t noise**: noise, outliers are not as unlikely as for Gaussian noise and thus will not have that much of an influence on the model.



Speed of Decay: $P(|y - \mathbf{w}^\top \mathbf{x}| > t)$ probability of having a outlier/derivation of larger than t , for linear regression.

Students-t $P(|y - \mathbf{w}^\top \mathbf{x}| > t) = \mathcal{O}(t^{-\alpha})$ ($\alpha > 0$) (Polynomial decay)

Gaussian $P(|y - \mathbf{w}^\top \mathbf{x}| > t) = \mathcal{O}(\exp^{-\alpha t})$ ($\alpha > 0$) (Exponential decay)

\Rightarrow **Students-t** distribution decays less fast than the Gaussian distribution and **thus** has heavier tails/tailmasses and does not get so easily influenced by noise.

Thus if we know that our model contains outliers/noise, we should use student's t distribution.

4. Proofs

Proof 4.2 4.12: From eq. (4.12) it follows that the response variables are uncorrelated given the explanatory variables $\text{Cov}[Y_i, Y_j | \mathbf{X}] = 0$. Hence we have i.i.d. samples with a corresponding conditional (log)-likelihood given by:

$$\begin{aligned} L_n(\mathbf{y}|\mathbf{X}, \theta) \stackrel{\text{i.i.d.}}{=} \prod_{i=1}^n P(\mathbf{x}_i, y_i | \theta) &= \prod_{i=1}^n \mathcal{N}(\beta^\top \mathbf{x}_i, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} \exp \left(-\frac{(y_i - \beta^\top \mathbf{x}_i)^2}{2\sigma^2} \right) \\ &= (\sigma^2 2\pi)^{-\frac{n}{2}} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^\top \mathbf{x}_i)^2 \right) \end{aligned}$$

$$l_n(\mathbf{y}|\mathbf{X}, \theta) = -\frac{n}{2} \ln \sigma^2 - \frac{n}{2} \ln 2\pi - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^\top \mathbf{x}_i)^2$$

Proof 4.3 Definition 4.14:

$$\begin{aligned} \beta^* \in \arg \min_{\beta \in \mathbb{R}^p} -l_n(\mathbf{y}|\mathbf{X}, \theta) \\ &= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^\top \mathbf{x}_i)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2\sigma^2} (\mathbf{y} - \beta^\top \mathbf{X})^\top (\mathbf{y} - \beta^\top \mathbf{X}) \\ &= \arg \min_{\beta \in \mathbb{R}^p} (\mathbf{y} - \beta^\top \mathbf{X})^\top (\mathbf{y} - \beta^\top \mathbf{X}) \\ &\stackrel{*}{\iff} (-2\mathbf{y}^\top \mathbf{X} + 2\mathbf{X}^\top \mathbf{X}\beta^*) = 0 \\ &\Rightarrow \mathbf{X}^\top \mathbf{X}\beta^* = \mathbf{X}^\top \mathbf{y} \end{aligned}$$

Note: *

$$\begin{aligned}
 & (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \\
 &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\beta + (\mathbf{X}\beta)^T \mathbf{y} - (\mathbf{X}\beta)^T (\mathbf{X}\beta) \\
 &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T (\mathbf{X}\beta) \\
 \frac{\partial}{\partial \mathbf{x}} \mathbf{Mx} = \mathbf{M} \quad \text{and} \quad \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{Mx} = (\mathbf{M} + \mathbf{M}^T)\mathbf{x} \quad (4.46)
 \end{aligned}$$

If we let $\mathbf{M} = \mathbf{X}^T \mathbf{X}$ then it follows:

$$\frac{\partial}{\partial \beta} \beta^T \mathbf{X}^T (\mathbf{X}\beta) = (\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T)\beta = 2\mathbf{X}^T \mathbf{X}\beta$$

Thus

$$0 = \frac{\partial}{\partial \beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) = 2\mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}) \quad (4.47)$$

Proof 4.4: [def. 4.13]

$$\begin{aligned}
 \text{lsq}(\mathbf{X}, \mathbf{y}) &= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \\
 &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\beta + (\mathbf{X}\beta)^T \mathbf{y} - (\mathbf{X}\beta)^T (\mathbf{X}\beta) \\
 &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T (\mathbf{X}\beta) \\
 0 &= \frac{\partial}{\partial \beta} \text{lsq}(\mathbf{X}, \mathbf{y}) = 2\mathbf{X}^T \mathbf{X}\beta - 2\mathbf{X}^T \mathbf{y} = 2\mathbf{X}^T (\mathbf{X}\beta - \mathbf{y})
 \end{aligned}$$

Note

$$\frac{\partial}{\partial \beta} \beta^T \mathbf{X}^T (\mathbf{X}\beta) \stackrel{\text{eq. (59.134)}}{=} (\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T)\beta = 2\mathbf{X}^T \mathbf{X}\beta$$

Proof 4.5: Corollary 4.2

$$\begin{aligned}
 & (\mathbf{X}\beta - \mathbf{y}) \perp \mathfrak{R}(\mathbf{X}) \\
 \iff & (\mathbf{X}\beta)^T (\mathbf{X}\beta - \mathbf{y}) = 0 \quad \forall \beta \in \mathbb{R}^m \\
 \iff & \mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}) = 0
 \end{aligned}$$

where $\mathbf{X} = \{\mathbf{x}_{:,1}, \dots, \mathbf{x}_{:,m}\}$ is the "basis" of the Range space:
 $(\mathbf{X}\beta - \mathbf{y})^T \mathbf{x}_{:,j} = 0 \quad \forall j = 1, \dots, m$

Proof 4.6 Corollary 4.3: From [def. 4.13] it follows:

$$\begin{aligned}
 \mathbf{X}^T \mathbf{Y} &= \mathbf{X}^T \mathbf{X}\hat{\beta} = \hat{\beta}^T \mathbf{X}^T \mathbf{X} = (\mathbf{X}\hat{\beta})^T \mathbf{X} \\
 \iff & (\mathbf{Y} - \mathbf{X}\hat{\beta})\mathbf{X} = \mathbf{r}^T \mathbf{X} = 0
 \end{aligned}$$

Proof 4.7 Property 4.4: $\hat{\beta}$ an unbiased estimator of β :

$$\begin{aligned}
 \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta + \epsilon) \\
 &= \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \\
 \mathbb{E}_\epsilon[\hat{\beta}] &= \mathbb{E}[\beta] + \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \\
 &= \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underbrace{\mathbb{E}[\epsilon]}_{=0} = \beta
 \end{aligned}$$

Proof 4.8 Property 4.4: Covariance $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$:

$$\begin{aligned}
 \text{Cov}[\hat{\beta}] &= \overbrace{\text{Cov}[\beta]}^{=0} + \overbrace{\text{Cov}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]}^{:= \mathbb{V}[\alpha\epsilon]} = \mathbb{E}[(\alpha\epsilon)^2] - \overbrace{\mathbb{E}[\alpha\epsilon]^2}^{=0} \\
 &= \mathbb{E}[(\alpha\epsilon)^T (\alpha\epsilon)] = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})] \\
 &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\epsilon \epsilon^T] \mathbf{X} (\mathbf{X}^T \mathbf{X}) \\
 &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}
 \end{aligned}$$

Proof 4.9 Property 4.5: $\hat{\mathbf{y}}$ an unbiased estimator of \mathbf{y} :

$$\mathbb{E}_\epsilon[\hat{\mathbf{y}}] = \mathbb{E}[\mathbf{X}\hat{\beta} + \epsilon] = \mathbf{X}\mathbb{E}[\hat{\beta}] + 0 \stackrel{\text{eq. (4.25)}}{=} \mathbf{X}\beta = \mathbb{E}[\mathbf{y}]$$

Proof 4.10 Theorem 4.3: $\hat{\beta}$ is a linear operator w.r.t. to \mathbf{y} :

$$\begin{aligned}
 \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} =: \mathbf{Cy} = \mathbf{C}(\mathbf{X}\beta) \\
 &= \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon =: \tilde{\mathbf{C}}\epsilon + \beta
 \end{aligned}$$

5. Examples

Example 4.1 Simple Linear Regression:

$$\begin{aligned}
 p &= 2 \quad \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}
 \end{aligned}$$

Example 4.2 Simple Linear Quadratic Regression:

$$\begin{aligned}
 p &= 3 \quad \mathbf{X} = \begin{pmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}
 \end{aligned}$$

Classification

6. Intro

Definition 4.19 Training Data

$$\mathcal{D} := \{(x_i, y_i) \mid x_i \in \mathcal{X} \subset \mathbb{R}^d, y_i \in \mathcal{Y} := \{c_1, \dots, c_K\}\}$$

Definition 4.20 Classifier

c:
Is a mapping that maps the features into classes:
 $c: \mathcal{X} \rightarrow \mathcal{Y}$

Definition 4.21 Dichotomy:

Given a set $\mathcal{S} = \{s_1, \dots, s_N\}$ a dichotomy is partition of the set \mathcal{S} into two subsets A, A^c that satisfy:

- Collectively/jointly exhaustiveness: $S = A \cup A^c$

- Mutual exclusivity: $s \in A \implies s \notin A^c \quad \forall s \in S$

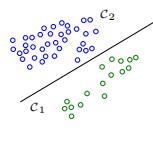
Explanation 4.1. Nothing can belong simultaneously to both parts A and A^c .

Types of Classification

Definition 4.22 Binary Classification:

Is a classification problem where the labels are binary:

$$\mathcal{Y} = \{c_1, c_2\} = \{-1, 1\} \quad (4.51)$$



Types of Categorical Data

Definition 4.23 Nominal/Categorical Data:

Is data where variables belong to a finite set of classes $\{c_1, \dots, c_K\}$ that do not have any ordering.

Definition 4.24 Ordinal Data:

Is data where variables belong to a finite discrete set of classes $\{c_1, \dots, c_K\}$ that are ordered/do have an ranking between each other i.e. numbers.

Encodings

6.3.1. Ordinal Encoding

Definition 4.25 Ordinal Encoding:

Each category gets assigned an integer values to introduce an order to the data.

Usage: for ordinal data, where we want to preserve order.

Cons

- models such as neural networks output a continues value, thus we are in fact treating a multil-class classification problem as regression problem.

6.3.2. One Hot Encoding

Definition 4.26 One-hot encoding/representation:

Is the representation/encoding of the K categories $\{c_1, \dots, c_K\}$ by a sparse vectors^[def. 59.70] with one non-zero entry, where the index j of the non-zero entry indicates the class c_j :

$$\mathbb{B}^n = \left\{ \mathbf{y} \in \{0, 1\}^n : \mathbf{y}^\top \mathbf{y} = \sum_{i=1}^n \mathbf{y}_i = 1 \right\}$$

s.t. $\mathbf{y}_i = \mathbf{e}_j \iff \mathbf{y}_i = \mathbf{c}_j$

Usage: for data where we do not want any order.

MNIST

I.e. for digit recognition we should treat our numbers as a set we do care that a 9 is classified as 9 but do not care that it comes after an .

6.3.3. Soft vs. Hard Labels

Definition 4.27 Hard Labels/Targets: Are observations $y \in \mathcal{Y}$ that are consider as true observations. We can encode them using a one hot encoding^[def. 4.26]:

$$y = \mathbf{c}_k \implies y = \mathbf{e}_k \quad (4.52)$$

Definition 4.28 Soft Labels/Targets: Are observations $y \in \mathcal{Y}$ that are consider as noisy observations or probabilities p. We can encode them using a probabilistic vector^[def. 59.71]: p. We can encode them using a probabilistic vector^[def. 59.71]:

$$y = [\mathbf{p}_1, \dots, \mathbf{p}_K]^\top \quad (4.53)$$

Corollary 4.8 Hard labels as special case: If we consider hard targets^[def. 4.27] as events with probability one then we can think of them as a special case of the soft labels.

7. Binary Classification

$$\{-1, 1\}$$

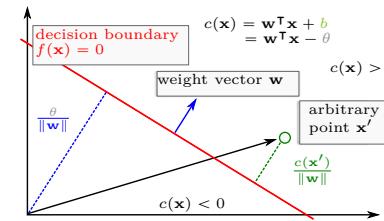
1. Linear Classification

Definition 4.29 Linear Dichotomy:

Definition 4.30 Linear Classifier: A linear classifier c that assigns labels \hat{y} to samples x_i using a linear decision boundary/hyperplane^[def. 59.15]:

$$\hat{y} = c(\mathbf{x}_i) = \begin{cases} c_1 \in \mathcal{H}^+ & \text{if } \mathbf{w}^\top \mathbf{x} > \theta \\ c_2 \in \mathcal{H}^- & \text{if } \mathbf{w}^\top \mathbf{x} < \theta \end{cases} \quad (4.54)$$

Explanation 4.2 (Definition 4.30):



- The $b \in \mathbb{R}$ corresponds to the offset of the decision surface from the origin, otherwise the decision surface would have to pass through the origin.
- $\mathbf{w} \in \mathbb{R}^d$ is the normal unit vector of the decision surface. Its components $\{w_j\}_{j=1}^d$ correspond to the importance of each feature/dimension.

Explanation 4.3 (Threshold θ vs. Bias b): The offset is called bias if it is considered as part of the classifier $\mathbf{w}^\top \mathbf{x} + b$ and as threshold if it is considered to be part of the hyperplane $\theta = -b$, but its just a matter of definition.

Definition 4.31 (Normalized) Classification Criterion:

$$\mathbf{w}^\top \mathbf{x} = \mathbf{w}^\top \mathbf{y} > 0 \quad \forall (\mathbf{x}, y) \in \mathcal{D} \quad (4.55)$$

Definition 4.32 Linear Separable Data set: A data set is linearly separable if there exists a separating hyperplane \mathcal{H} s.t. each label can be assigned correctly:

$$\hat{y} := c(\mathbf{x}) = y \quad \forall (\mathbf{x}, y) \in \mathcal{D} \quad (4.56)$$

7.1.1. Normalization

Proposition 4.3 Including the Offset: In order to simplify notation the offset is usually included into the parameter vector:

$$\mathbf{w} \leftarrow \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \quad \mathbf{x} \leftarrow \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}$$

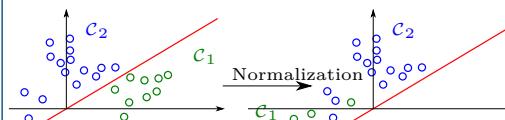
$$\Rightarrow \mathbf{w}^\top \mathbf{x} = (\mathbf{w}^\top \mathbf{x}) \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} = \mathbf{w}^\top \mathbf{x} + b$$

Proposition 4.4 Uniform Classification Criterion: In order to avoid the case distinction in the classification criterion of eq. (4.54) we may transform the input samples by:

$$\tilde{\mathbf{x}} = \begin{cases} \mathbf{x} & \text{if } \mathbf{w}^\top \mathbf{x} > \theta \\ -\mathbf{x} & \text{if } \mathbf{w}^\top \mathbf{x} < \theta \end{cases} \quad (4.57)$$

Explanation 4.4 (proposition 4.4).

We transform the input s.t. the separating hyper-plane puts all labels on the same “positive” side $\mathbf{w}^\top \mathbf{x} > 0$.



Corollary 4.9 : How can we achieve this in practice?

If $\mathcal{Y} = \{-1, 1\}$ then we can simply multiply with the label y_i :

$$\begin{cases} \mathbf{w}^\top \mathbf{x} > 0 & \forall y = +1 \\ \mathbf{w}^\top \mathbf{x} < 0 & \forall y = -1 \end{cases} \iff \mathbf{w}^\top \mathbf{x} \cdot y > 0 \quad \forall y$$

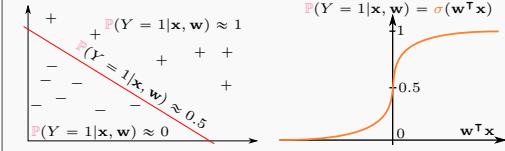
8. Logistic Regression

$$\text{Bern}(y; \sigma(\mathbf{w}^\top \mathbf{x}, \sigma^2))$$

Idea: in order to classify dichotomies^[def. 4.21] we use a distribution that maps probabilities to a binary values 0/1 \Rightarrow Bernoulli Distribution^[def. 66.22].

Problem: we need to convert/translate distance $\mathbf{w}^\top \mathbf{x}$ into probability in order to use a bernoulli distribution.

Idea: use a sigmoidal function to convert distances $z := \mathbf{w}^\top \mathbf{x}$ into probabilities \Rightarrow Logistic Function^[def. 4.33].



1. Logistic Function

Definition 4.33 Sigmoid/Logistic Function:

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-\text{neg. dist. from deci. boundary}}} \quad (4.58)$$

Explanation 4.5 (Sigmoid/Logistic Function).

$$\sigma(z) = \begin{cases} 0 & \text{if } z \text{ large} \\ 1 & \text{if } z \text{ large} \\ 0.5 & z = 0 \end{cases}$$

2. Logistic Regression

Definition 4.34 Logistic Regression: models the likelihood of the output y as a Bernoulli Distribution^[def. 66.22] $y \sim \text{Bern}(p)$, where the probability p is given by the Sigmoid function^[def. 4.33] of a linear regression:

$$\begin{aligned} P(y|\mathbf{x}, \mathbf{w}) &= \text{Bern}(\sigma(\mathbf{w}^\top \mathbf{x})) = \begin{cases} \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} & \text{if } y = +1 \\ \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} & \text{if } y = -1 \end{cases} \\ &\stackrel{??}{=} \frac{1}{1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x})} = \sigma(-y \cdot \mathbf{w}^\top \mathbf{x}) \end{aligned} \quad (4.59)$$

8.2.1. Maximum Likelihood Estimate

Definition 4.35 Logistic Loss l_L proof 4.12:
Is the objective we want to minimize when performing mle^[def. 6.3] for a logistic regression likelihood and incurs higher cost for samples closer to the decision boundary:

$$l_L(\mathbf{w}; \mathbf{x}, y) := \log(1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x})) \quad (4.60)$$

$$\propto \log(1 + e^z) = \begin{cases} z & \text{for large } z \\ 0 & \text{for small } z \end{cases}$$

Corollary 4.10 MLE for Logistic Regression:

$$l_n(\mathbf{w}) = \sum_{i=1}^n l_i = \sum_{i=1}^n \log(1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)) \quad (4.61)$$

Stochastic Gradient Descent

The logistic loss l_L is a convex function. Thus we can use convex optimization techniques s.a. SGD in order to minimize the objective [cor. 4.10].

Definition 4.36

Logistic Loss Gradient

$$\nabla_{\mathbf{w}} l_i(\mathbf{w}) = \nabla_{\mathbf{w}} \log(1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)) = \frac{1}{1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)} \cdot (-\mathbf{x}_i) \quad (4.62)$$

Explanation 4.6.

$$\nabla_{\mathbf{w}} l_i(\mathbf{w}) = \nabla_{\mathbf{w}} \log(1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i))$$

The logistic loss l_L is equal to the hinge loss l_H but weighted by the probability of being in the wrong class $P(Y = -1|\mathbf{x}, \mathbf{w})$. Thus the more likely we are in the wrong class the bigger the step we take:

$$P(Y = -1|\hat{y} = \mathbf{w}^\top \mathbf{x}) = \begin{cases} \uparrow & \text{take big step} \\ \downarrow & \text{take small step} \end{cases}$$

Algorithm 4.1 Vanilla SGD for Logistic Regression:

Initialize: \mathbf{w}

- 1: for $i = 1, 2, \dots, T$ do
- 2: Pick (\mathbf{x}_i, y_i) unif. at random from data \mathcal{D}
- 3: $\mathbf{P}(Y = -1|\mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)} = \sigma(y_i \cdot \mathbf{w}^\top \mathbf{x}_i)$ compute prob. of misclassif. with cur. model
- 4: $\mathbf{w} = \mathbf{w} + \eta_{t,y} \mathbf{x}_i \sigma(y_i \cdot \mathbf{w}^\top \mathbf{x}_i)$
- 5: end for

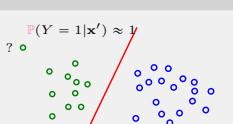
Making Predictions

Given an optimal parameter vector $\hat{\mathbf{w}}$ found by algorithm 4.1 we can predict the output of a new label by eq. (4.59):

$$P(y|\mathbf{x}, \hat{\mathbf{w}}) = \frac{1}{1 + \exp(-y \cdot \hat{\mathbf{w}}^\top \mathbf{x})} \quad (4.63)$$

Drawback

Logistic regression, does not tell us anything about the likelihood $P(\mathbf{x})$ of a point, thus it will not be able to detect outliers, as it will assign a very high probability to all correctly classified points, far from the decision boundary.



8.2.2. Maximum a-Posteriori Estimates

3. Logistic regression and regularization

Adding Priors to Logistic Likelihood

- **L2 (Gaussian prior):** arg min $\sum_{i=1}^n \log(1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)) + \lambda \|\mathbf{w}\|_2^2$
- **L1 (Laplace prior):** arg min $\sum_{i=1}^n \log(1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)) + \lambda \|\mathbf{w}\|_1$
- **Generalized:** $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{i=1}^n \log(1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)) + \lambda C(\mathbf{w})$
 $= \arg \max_{\mathbf{w}} P(\mathbf{w}|\mathbf{X}, Y)$

4. SGD for L2-regularized logistic regression

```

Initialize:  $\mathbf{w}$ 
1: for  $1, 2, \dots, T$  do
2:   Pick  $(\mathbf{x}, y)$  unif. at random from data  $\mathcal{D}$ 
3:    $\hat{P}(Y = -y | \mathbf{x}, \mathbf{w}) = \frac{1}{(1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x}))}$ 
       $\triangleright$  compute prob. of misclassif. with cur. model
4:    $\mathbf{w} = \mathbf{w} (1 - 2\lambda \eta_t) + \eta_t y \mathbf{x} \hat{P}(Y = -y | \mathbf{x}, \mathbf{w})$ 
5: end for

```

Thus: \mathbf{w} is pulled/shrunken towards zero, depending on the regularization parameter $\lambda > 0$

9. Proofs

Proof 4.11: [def. 4.34] We need to only proof the second expression, as the first one is fulfilled anyway:

$$1 - \frac{1}{1 + e^z} = \frac{1 + e^z}{1 + e^z} - \frac{1}{1 + e^z} = \frac{e^z + 1 - 1}{1 + e^z} = \frac{e^z}{e^z + 1}$$

$$= \frac{1}{1 + e^{-z}}$$

Proof 4.12: [def. 4.35]

$$\begin{aligned} l_n(\mathbf{w}) &= \arg \max_{\mathbf{w}} P(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) = \arg \min_{\mathbf{w}} -\log P(Y | \mathbf{X}, \mathbf{w}) \\ &\stackrel{\text{i.i.d.}}{=} \arg \min_{\mathbf{w}} \sum_{i=1}^n -\log P(y_i | \mathbf{x}_i, \mathbf{w}) \\ &\stackrel{\text{eq. (4.59)}}{=} -\log \frac{1}{1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)} \\ &= \log(1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)) =: l_i(\mathbf{w}) \end{aligned}$$

Proof 4.13: [def. 4.36]

$$\begin{aligned} \nabla_{\mathbf{w}} l_i(\mathbf{w}) &= \frac{\partial}{\partial \mathbf{w}} \log(1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x})) \\ &\stackrel{\text{C.R.}}{=} \frac{1}{(1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x}))} \frac{\partial}{\partial \mathbf{w}} (1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x})) \\ &\stackrel{\text{C.R.}}{=} \frac{1}{(1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x}))} \exp(-y \cdot \mathbf{w}^\top \mathbf{x}) \cdot (-y \mathbf{x}) \\ &= \frac{e^{-z} \cdot (-yx)}{(1 + e^{-z})} = \frac{-yx}{e^z(1 + e^{-z})} = \frac{-yx}{(e^z + e^{-z+z})} \\ &= \frac{1}{\exp(y \cdot \mathbf{w}^\top \mathbf{x}) + 1} \cdot (-yx) \\ &\stackrel{\text{eq. (4.59)}}{=} \hat{P}(Y = -y | \mathbf{x}, \mathbf{w}) \cdot (-yx) \end{aligned}$$

Generalized Linear Models (GLMs)

1. Generalized Additive Models (GAMs)

Definition 5.1 $g_{\text{add}} : \mathbb{R}^p \mapsto \mathbb{R}$

Generalized Additive Models (GAM):
Are generalized linear model where the response variable depends linearly on unknown smooth functions g_j s.t.:

$$g_{\text{add}}(\mathbf{x}) = \mu + \sum_{j=1}^p g_j(x_j) \quad g_j : \mathbb{R} \mapsto \mathbb{R} \quad \forall j \in \{1, \dots, p\}$$
$$\mathbb{E}[g_j(x_j)] = 0 \quad (5.1)$$

Pros

- Does not suffer from the curse of dimensionality.

Cons

- does not allow for interaction terms such as $g_{j,k}(x_j, x_k)$.

1. Backfitting

Model Parameter Estimation

1. Maximum Likelihood Estimation

1. Likelihood Function

Is a method for estimating the parameters θ of a model that agree best with observed data $\{x_1, \dots, x_n\}$. Let: $\theta = (\theta_1 \dots \theta_k)^\top \in \Theta \subset \mathbb{R}^k$ vector of unknown model parameters.

Consider: a probability density/mass function $f_X(\mathbf{x}; \theta)$

Definition 6.1 Likelihood Function $L_n : \Theta \times \mathbb{R}^n \mapsto \mathbb{R}_+$: Let $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ be a random sample of i.i.d. data points drawn from an unknown probability distribution $\mathbf{x}_i \sim p_{\mathcal{X}}$. The likelihood function gives the likelihood/probability of the joint probability of the data $\{x_1, \dots, x_n\}$ given a fixed set of model parameters θ :

$$L_n(\theta | \mathbf{X}) = L_n(\theta; \mathbf{X}) = f(\mathbf{X} | \theta) = f(\mathbf{X}; \theta) \quad (6.1)$$

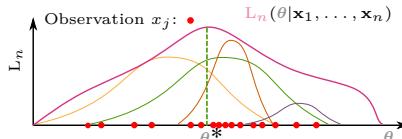


Figure 5: Possible Likelihood function in pink.
Overlayed: possible candidate functions for Gaussian model explaining the observations.

Likelihood function is not a pdf

The likelihood function by default not a probability density function and may not even be differentiable. However if it is, then it may be normalized to one.

Corollary 6.1 i.i.d. data: If the n-data points of our sample are i.i.d. then the likelihood function can be decomposed into a product of n-terms:

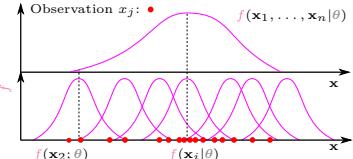


Figure 6: Bottom: probability distributions of the different data points \mathbf{x}_i given a fixed θ for a Gaussian distribution
Top: joint probability distribution of the i.i.d. data points $\{\mathbf{x}_i\}_{i=1}^n$ given a fixed θ

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n | \theta) \stackrel{i.i.d.}{=} \prod_{i=1}^n f(\mathbf{x}_i | \theta)$$

Notation

- The probability density $f(\mathbf{X} | \theta)$ is considered for a fixed θ and thus as a function of the samples.
- The likelihood function on the other hand is considered as a function over parameter values θ for a fixed sample $\{\mathbf{x}_i\}_{i=1}^n$ and thus written as $L_n(\theta | \mathbf{X})$.
- Often the colon symbol ; is written instead of the is given symbol | in order to indicate that θ resp. \mathbf{X} is a parameter and not a random variable.

2. Maximum Likelihood Estimation (MLE)

Let $f_\theta(\mathbf{x})$ be the probability of an i.i.d. sample \mathbf{x} for a given model.

Goal: find θ of a given model that maximizes the joint probability/likelihood of the observed data $\{x_1, \dots, x_n\}$? \iff maximum likelihood estimator θ^* .

Definition 6.2 Log Likelihood Function

$$l_n(\theta | \mathbf{X}) = \log L_n(\theta | \mathbf{X}) = \log f(\mathbf{X} | \theta) \quad (6.2)$$

Corollary 6.2 i.i.d. data: Differentiating the product of n-Terms with the help of the chain rule leads often to complex terms. As a result one usually prefers maximizing the log (especially for exponential terms), as it does not change the argmax—eq. (54.66):

$$\log f(\mathbf{x}_1, \dots, \mathbf{x}_n | \theta) \stackrel{i.i.d.}{=} \log \left(\prod_{i=1}^n f(\mathbf{x}_i | \theta) \right) = \sum_{i=1}^n \log f(\mathbf{x}_i | \theta)$$

Definition 6.3 Maximum Likelihood Estimator

θ^* : Is the estimator $\theta^* \in \Theta$ that maximizes the likelihood of the model/predictor:

$$\theta^* = \arg \max_{\theta \in \Theta} L_n(\theta; \mathbf{X}) \quad \text{or} \quad \theta^* = \arg \max_{\theta \in \Theta} l_n(\theta; \mathbf{X}) \quad (6.3)$$

3. Maximization vs. Minimization

For optimization problems we minimize by convention. The logarithm is a concave function^[def. 54.25] \cap , thus if we calculate the extremal point we will obtain a maximum. If we want to calculate a minimum instead (i.e. in order to be compatible with some computer algorithm) we can convert the function into a convex function^{section 5} \cup by multiplying it by minus one and consider it as a loss function instead of a likelihood.

Definition 6.4 Negative Log-likelihood

$$-\bar{l}_n(\theta | \mathbf{X}): \theta^* = \arg \max_{\theta \in \Theta} l_n(\theta | \mathbf{X}) = \arg \min_{\theta \in \Theta} -l_n(\theta | \mathbf{X}) \quad (6.4)$$

4. Conditional Maximum Likelihood Estimation

Maximum likelihood estimation can also be used for conditional distributions.

Assume the labels y_i are drawn i.i.d. from a unknown true conditional probability distribution $f_{Y|X}$ and we are given a data set $\mathbf{Z} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}\}_{i=1}^n$.

Now we want to find the parameters $\theta = (\theta_1 \dots \theta_k)^\top \in \Theta \subset \mathbb{R}^k$ of a hypothesis $\hat{f}_{Y|X}$ that agree best with the given data \mathbf{Z} .

Note

For simplicity we omit the hat $\hat{\cdot}$ of our model $\hat{f}_{Y|X}$ and simply assume that our data is generated by some data generating probability distribution.

Definition 6.5 Conditional (log) likelihood function:

Models the liklihood of a model with parameters θ given the data $\mathbf{Z} = \{\mathbf{x}_i, y_i\}_{i=1}^n$

$$L_n(\theta | Y, \mathbf{X}) = L_n(\theta; Y, \mathbf{X}) = f(Y | \mathbf{X}, \theta) = f(Y | \mathbf{X}; \theta)$$

2. Maximum a posteriori estimation (MAP)

Idea

We have seen (??), that trading/increasing a bit of bias can lead to a big reduction of variance of the generalization error. We also know that the least squares MLE is unbiased (??). Thus the question arises if we can introduce a bit of bias into the MLE in turn of decreasing the variance?
 \Rightarrow use Bayes rule (??) to introduce a bias into our model via. a Prior distribution.

1. Prior Distribution

Definition 6.6 Prior (Distribution)

$\pi(\theta) = p(\theta)$: Assumes: that the model parameters θ are no longer constant but random variables distributed according to a prior distribution that models some prior belief/bias that we have about the model:

$$\theta \sim \pi(\theta) = p(\theta) \quad (6.5)$$

Notes

In this section we use the terms model parameters θ and model as synonymous, as the model is fully described by its population parameters^[def. 3.12] θ .

Corollary 6.3 The prior is independent of the data:

The prior $p(\theta)$ models a prior belief/bias and is thus independent of the data $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$:

$$p(\theta | \mathbf{X}) = p(\theta) \quad (6.6)$$

Definition 6.7 Hyperparameters

$p_{\lambda}(\theta)$: In most cases the prior distribution are parameterized that is the pdf $p(\theta | \lambda)$ depends on a set of parameters λ . The parameters of the prior distribution, are called hyperparameters and are supplied due to believe/prior knowledge (and do not depend on the data) see example 6.1

2. Posterior Distribution

Definition 6.8 Posterior Distribution

$p(\theta | \text{data})$: The posterior distribution $p(\theta | \text{data})$ is a probability distribution that describes the relationship of a unknown parameter θ a posterior/after observing evidence of a random quantity Z that is in a relation with θ :

$$p(\theta | \text{data}) = p(\theta | Z) \quad (6.7)$$

Definition 6.9 [proof 22.1]

Posterior Distribution and Bayes Theorem:

Using Bayes theorem 65.3 we can write the posterior distribution as a product of the likelihood^[def. 6.1] weighted with our prior^[def. 6.6] and normalized by the evidence $Z = \{\mathbf{X}, \mathbf{y}\}$ s.t. we obtain a real probability distribution:

$$p(\theta | \text{data}) = p(\theta | Z) = \frac{p(Z | \theta) \cdot p_{\lambda}(\theta)}{p(Z)} \quad (6.8)$$

$$\text{Posterior} = \frac{\text{Liklihood} \cdot \text{Prior}}{\text{Normalization}} \quad (6.9)$$

$$p(\theta | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \theta, \mathbf{X}) \cdot p_{\lambda}(\theta)}{p(\mathbf{y} | \mathbf{X})} \quad (6.10)$$

2.2.1. Maximization –MAP

We do not care about the full posterior probability distribution as in Bayesian Inference (section 3). We only want to find a point estimator ?? θ^* that maximizes the posterior distribution.

2.2.2. Maximization

Definition 6.10 Maximum a-Posteriori Estimates (MAP):

Is model/parameters θ that maximize the posterior probability distribution:

$$\theta_{\text{MAP}}^* = \arg \max_{\theta} p(\theta | \mathbf{X}, \mathbf{y}) \quad (6.11)$$

Log-MAP estimator:

$$\theta^* = \arg \max_{\theta} \{p(\theta | \mathbf{X}, \mathbf{y})\} \quad (6.12)$$

$$= \arg \max_{\theta} \left\{ \frac{p(\mathbf{y} | \mathbf{X}, \theta) \cdot p_{\lambda}(\theta)}{p(\mathbf{y} | \mathbf{X})} \right\}$$

$$\underset{\text{eq. (54.63)}}{\approx} \arg \max_{\theta} \left\{ p(\mathbf{y} | \theta, \mathbf{X}) \cdot p_{\lambda}(\theta) \right\}$$

Corollary 6.4 Negative Log MAP:

$$\theta^* = \arg \max_{\theta} \{p(\theta | \mathbf{X}, \mathbf{y})\} \quad (6.13)$$

$$= \arg \min_{\theta} -\log \overbrace{p(\theta)}^{\text{Prior}} - \log \overbrace{p(\mathbf{y} | \theta, \mathbf{X})}^{\text{Likelihood}} + \underbrace{\log p(\mathbf{y} | \mathbf{X})}_{\text{not depending on } \theta}$$

3. Examples

Example 6.1 Hyperparameters Gaussian Prior:

$$f_{\lambda}(\theta) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{(\theta - \mu)^2}{2\sigma^2} \right)$$

with the hyperparameter $\lambda = (\mu \ \sigma^2)^\top$.

Dimensionality Reduction

Bayesian Inference/Modeling

Definition 6.11 Bayesian Inference: So far we only really looked at point estimators/estimates^[def. 6.8.8]. But what if we are interested not only into the most likely value but also want to have a notion of the uncertainty of our prediction? Bayesian inference refers to statistical inference^[def. 3.10], where uncertainty in inferences is quantified using probability. Thus we usually obtain a distribution over our parameters and not a single point estimates \Rightarrow can deduce statistical properties of parameters from their distributions.

Definition 6.12 $p(w|y, X)/p(w|\mathcal{D})$

Posterior Probability Distribution:

- ① Specify the prior $p_A(w)$
 - ② Specify the likelihood $p(y|w, X)/p(\mathcal{D}|w)$
 - ③ Calculate the evidence $p(y|X)/p(\mathcal{D})$
 - ④ Calculate the posterior distribution $P(w|y, X)/P(w|\mathcal{D})$
- $$p(w|y, X) = \frac{p(y|w, X) \cdot p_A(w)}{p(y|X)} = \frac{\text{Likelihood} \cdot \text{Prior}}{\text{Normalization}}$$

Definition 6.13 $p(y|X)/p(\mathcal{D})$

Marginal Likelihood [see proof 10.2]: is the normalization constant that makes sure that the posterior distribution^[def. 6.12] is a true probability distribution:

$$p(y|X) = \int p(y|w, X) \cdot p_A(w) dw = \int \text{Likelihood} \cdot \text{Prior} dw \quad (6.14)$$

Note

It is called marginal likelihood as we marginalize over w .

Definition 6.14 Posterior Marginal Distribution: Is the posterior distribution of single elements of our thought after parameter vector:

$$p(w_i|y, X) = \int p(y|w, X) dw_{-i} \quad i = 1, \dots, \dim(w) \quad (6.15)$$

Definition 6.15 $p(f_*|x_*, X, y)/p(f_*|y)$ [see proof 10.1]

Posterior Predictive Distribution:

is the distribution of a real process f (i.e. $f(x) = x^T w$) given:

- new observation(s) x_*
- the posterior distribution^[def. 6.12] of the observed data $\mathcal{D} = \{X, y\}$
- The likelihood of a real process f_*

$$p(f_*|x_*, X, y) = \int p(f_*|x_*, w) \cdot p(w|X, y) dw \quad (6.16)$$

it is calculated by weighting the likelihood^[def. 6.1] of the new observation x_* with the posterior of the observed data and averaging over all parameter values w . \Rightarrow obtain a distribution not depending on w .

Note f vs. y

- Usually f denotes the model i.e.: $f(x) = x^T w$ or $f(x) = \phi(x)^T w$ and y the model plus the noise $y = f(x) + \epsilon$.
- Sometimes people also write only: $p(y_*|x_*, X, y)$

4. Types of Uncertainty

Definition 6.16 Epistemic/Systematic Uncertainty:

Is the uncertainty that is due to things that one could in principle know but does not i.e. only having a finite sub sample of the data. The epistemic noise will decrease the more data we have.

Definition 6.17 Aleatoric/Statistical Uncertainty:

Is the uncertainty of an underlying random process/model. The aleatoric uncertainty stems from the fact that we are create random process models. If we run our *trained* model multiple times with *the same* input X data we will end up with different outcomes \hat{y} . The aleatoric noise is *irreducible* as it is an underlying part of probabilistic models.

Bayesian Filtering

Definition 7.1

Recursive Bayesian Estimation/Filtering: Is a technique for estimating the an unknown probability distribution recursively over time by a measurement^[def. 7.3] and a process-model^[def. 7.2] using Bayesian inference^[def. 6.11].

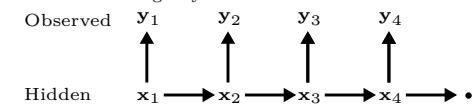


Figure 7: This problem corresponds to a *hidden Markov model (HMM)*^[def. 14.1]

$$\mathbf{x}_t = (x_{t,1} \dots x_{t,n}) \quad \mathbf{y}_t = (y_{t,1} \dots y_{t,m})$$

Note

Comes from the idea that spam can be filtered out by the probability of certain words.

Definition 7.2

$x_{t+1} \sim p(x_t|x_{t-1})$
Process/Motion/Dynamic Model: is a model q of how our system state x_t evolves and is usually fraught with some uncertainty.

Corollary 7.1 Markov Property $x_t \perp\!\!\!\perp \mathbf{x}_{1:t-1} | x_{t-1}$: The process models^[def. 7.2] is Markovian^[def. 69.14] i.e. the current state depends only on the previous state:

$$p(x_t|x_{1:t-1}) = p(x_t|x_{t-1}) \quad (7.1)$$

Definition 7.3

$$y_t \sim p(y_t|x_t)$$

Measurement/Sensor-Model/Likelihood: is a model h that maps observations/sensor measurements of our model \mathbf{y}_t to the model state x_t

Corollary 7.2

$y_t \perp\!\!\!\perp \mathbf{x}_{1:t-1} | x_t$
Conditional Independent Measurements: The measurements y_t are conditionally independent of the previous observations $\mathbf{x}_{1:t-1}$ given the current state x_t :

$$p(y_t|\mathbf{x}_{1:t-1}, x_t) = p(y_t|x_t) \quad (7.2)$$

Goal

We want to combine the process model^[def. 7.2] and the measurement model^[def. 7.3] in a recursive way to obtain a good estimate of our model state:

$$p(x_t|x_{t-1}) \xrightarrow{\text{p(x_t|y_{1:t})}} p(x_{t+1}|y_{1:t}) \xrightarrow{\text{recursion rule}} p(x_{t+1}|y_{1:t+1})$$

Definition 7.4 Chapman-Kolmogorov eq. $p(x_t|y_{1:t-1})$

Prior Update/Prediction Step [see proof 10.3]:

$$p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1}) dx_{t-1} \quad (7.3)$$

Prior Distribution:

$$p(x_0|y_{0-1}) = p(x_0) = p_0 \quad (7.4)$$

Definition 7.5

$p(x_t|y_{1:t})$
Posterior Distribution/Update Step [see proof 10.4]:

$$p(x_t|y_{1:t}) = \frac{1}{Z_t} p(y_t|x_t) p(x_t|y_{1:t-1}) \quad (7.5)$$

Definition 7.6 Normalization [see proof 10.5]:

$$Z_t = p(y_t|y_{1:t-1}) = \int p(y_t|x_t) p(x_t|y_{1:t-1}) dx_t \quad (7.6)$$

Algorithm 7.1 Optimal Bayesian Filtering:

```

1: Input:  $p(x_0)$ 
2: while Stopping Criterion not full-filled do
3:   Prediction Step:
         $p(x_t|y_{1:t}) = \frac{1}{Z_t} p(y_t|x_t) p(x_t|y_{1:t-1})$ 
5:   Update Step:
         $p(x_t|y_{1:t-1}) = \int p(x_t|x_{t-1}) p(x_{t-1}|y_{1:t-1}) dx_{t-1}$ 
        with:
         $Z_t = \int p(y_t|x_t) p(x_t|y_{1:t-1}) dx_t$ 
6: end while

```

Corollary 7.3 [proof 10.6]

Joint Probability Distribution of (HMM): we can also calculate the joint probability distribution of the (HMM):

$$p(x_{1:t}, y_{1:t}) = p(x_1)p(y_1|x_1) \prod_{i=2}^t p(x_i|x_{i-1})p(y_i|x_i) \quad (7.7)$$

Example 7.1 Types of Bayesian Filtering:

- **Kalman Filter:** assumes a *linear* system, q, h are linear and Gaussian noise v, w .
- **Extended Kalman Filter:** assumes a *non-linear* system, q, h are non-linear and Gaussian noise v, w .
- **Particle Filter:** assumes a *non-linear* system q, h are non-linear and Non-Gaussian noise v, w , especially multi-modal distributions.

1. Kalman Filters

Definition 7.7 Kalman Filter Assumptions: Assumes a *linear*^[def. 54.15] process model^[def. 7.2], q with Gaussian model-noise v and a linear measurement model^[def. 7.3] h with Gaussian process-noise w .

Definition 7.8 Kalman Filter Model:

Process Model (7.8)

$$x[k] = A[k-1]x[k-1] + u[k-1] + v[k-1] \quad \text{with} \\ x[0] \sim \mathcal{N}(x_0, P_0) \quad \text{and} \quad v^{(k)} \sim \mathcal{N}(0, Q^{(k)})$$

Measurement Model (7.9)

$$z[k] = H[k]x[k] + w[k] \quad \text{with} \quad w[k] \sim \mathcal{N}(0, R^{(k)})$$

and define:

$$\hat{x}_p^{(k)} := \mathbb{E}[x_p^{(k)}] \quad \text{and} \quad P_p^{(k)} := \mathbb{V}[x_p^{(k)}] \quad (7.10)$$

$$\hat{x}_m^{(k)} := \mathbb{E}[x_m^{(k)}] \quad \text{and} \quad P_m^{(k)} := \mathbb{V}[x_m^{(k)}] \quad (7.11)$$

Note

The CRVs $x_0, \{v(\cdot)\}, \{w(\cdot)\}$ are mutually independent.

Gaussian Processes (GP)

1. Gaussian Process Regression

1. Gaussian Linear Regression

Given

(1) Linear Model with Gaussian Noise:
 $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$
 $\epsilon \sim \mathcal{N}(0, \sigma_n^2 \mathbf{I})$

$\mathbf{y} = f(\mathbf{x}) + \epsilon$

\Rightarrow Gaussian Likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma_n^2 \mathbf{I})$

(2) Gaussian Prior: $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_p)$

Sought

(1) Posterior Distribution: $p(\mathbf{w}|\mathbf{y}, \mathbf{X})$

(2) Posterior Predictive Distribution: $p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y})$

Definition 8.1 $p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \mathcal{N}(\mu_w, \Sigma_w)$

Posterior Distribution proof 10.7:
 $\mu_w = \frac{1}{\sigma_n^2} \Sigma_w^{-1} \mathbf{X} \mathbf{y}$
 $\Sigma_w = \frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}^\top + \Sigma_p^{-1}$

Note

We could also use a prior with non-zero mean $p(\mathbf{w}) = \mathcal{N}(\mu, \Sigma_p)$ but by convention w.o.l.g. we use zero mean see ??.

Definition 8.2 $p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mu_*, \Sigma_*)$

Posterior Predictive Distribution proof 10.8:
 $\mu_* = \frac{1}{\sigma_n^2} \mathbf{x}_*^\top \Sigma_w^{-1} \mathbf{X} \mathbf{y}$
 $\Sigma_* = \mathbf{x}_*^\top \Sigma_w^{-1} \mathbf{x}_*$

2. Kernelized Gaussian Linear Regression

Definition 8.3 Posterior Predictive Distribution:
 $p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mu_*, \Sigma_*)$

μ_* (8.4)

Definition 8.4 Gaussian Process:

2. Model Selection

1. Marginal Likelihood

Approximate Inference

Problem

In statistical inference we often want to calculate integrals of probability distributions i.e.

- Expectations

$$\mathbb{E}_{\mathbf{X} \sim p} [g(\mathbf{X})] = \int g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

- Normalization constants:

$$p(\theta|y) = \frac{1}{Z} p(\theta, y) = \frac{p(y|\theta)p(\theta)}{Z} = \frac{p(y|\theta)p(\theta)}{\int p(\theta) d\theta}$$

$$Z = \int p(\theta|y)p(\theta) d\theta = \int p(\theta) \prod_{i=1}^n p(y_i|\mathbf{x}_i, \theta) d\theta$$

For non-linear distributions this integrals are in general intractable which may be due to the fact that there exist no analytic form of the distribution we want to integrate or highly dimensional latent spaces that prohibits numerical integration (curse of dimensionality).

Definition 9.1 Approximate Inference: Is the procedure of finding an probability distribution q that approximates a true probability distribution p as well as possible.

1. Variational Inference

Definition 9.2 Bayes Variational Inference:

Given an unnormalized (posterior) probability distribution:

$$p(\theta|y) = \frac{1}{Z} p(\theta, y) \quad (9.1)$$

BVI seeks an approximate probability distribution q_λ , that is parameterized by a variational parameter λ and approximates $p(\theta|y)$ well.

Definition 9.3 Variational Family of Distributions Q : a set of probability distributions Q that is parameterized by the same variational parameter λ is called a variational family.

1. Laplace Approximation

Definition 9.4 [example 10.1], [proof 10.9, 10.10, 10.11]
Laplace Approximation: Tries to approximate a desired probability distribution $p(\theta|\mathcal{D})$ by a Gaussian probability distribution:

$$Q = \{q_\lambda(\theta) = \mathcal{N}(\lambda)\} = \mathcal{N}(\mu, \Sigma) \quad (9.2)$$

the distribution is given by:

$$q(\theta) = c \cdot \mathcal{N}(\theta; \lambda_1, \lambda_2) \quad (9.3)$$

$$\lambda_1 = \hat{\theta} = \arg \max_{\theta} p(\theta|y)$$

with

$$\lambda_2 = \Sigma = H^{-1}(\hat{\theta}) = -\nabla \nabla_\theta \log p(\hat{\theta}|y)$$

Note

The name *Laplace Approximation* comes from its inventor *Pierre-Simon Laplace*.

Corollary 9.1 : Taylor approximation of a function $p(\theta|y) \in \mathcal{C}^k$ around its mode $\hat{\theta}$ naturally induces a Gaussian approximation. See proofs 10.9, 10.10, 10.11

2. Black Box Stochastic Variational Inference

The most common way of finding q_λ is by minimizing the KL-divergence^[def. 3.8] between our approximate distribution q and our true posterior p :

$$q^* \in \arg \min_{q \in Q} \text{KL}(q(\theta) \parallel p(\theta|y)) = \arg \min_{\lambda \in \mathbb{R}^d} \text{KL}(q_\lambda(\theta) \parallel p(\theta|y))$$

Note

Usually we want to minimize $\text{KL}(p(\theta|y) \parallel q(\theta))$ but this is often infeasible s.t. we only minimize $\text{KL}(q(\theta) \parallel p(\theta|y))$

Definition 9.5

[proof 10.12]

ELBO-Optimization Problem:

$$q^* \in \arg \min_{\lambda: q_\lambda \in Q} \text{KL}(q_\lambda(\theta) \parallel p(\theta|y))$$

$$= \arg \max_{\lambda: q_\lambda \in Q} \mathbb{E}_{\theta \sim q_\lambda} [\log p(y, \theta)] + H(q_\lambda) \quad (9.4)$$

$$= \arg \max_{\lambda: q_\lambda \in Q} \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] - \text{KL}(q_\lambda(\theta) \parallel p(\theta)) \quad (9.5)$$

$$:= \arg \max_{\lambda: q_\lambda \in Q} \text{ELBO}(\lambda) \quad (9.6)$$

Attention: Sometimes people write simply p for the posterior and $p(\cdot)$ for prior.

Explanation 9.1.

- eq. (9.4):
 - prefer uncertain approximations i.e. we maximize $H(q)$
 - that jointly make the joint posterior likely
- eq. (9.6): Expected likelihood of our posterior over q minus a regularization term that makes sure that we are not too far away from the prior.

3. Expected Lower Bound of Evidence (ELBO)

Definition 9.6

[example 10.2]/[proof 10.13]

Expected Lower Bound of Evidence (ELBO):

The evidence lower bound is a bound on the log prior:

$$\text{ELBO}(q_\lambda) \leq \log p(y) \quad (9.7)$$

1.3.1. Maximizing The ELBO

Definition 9.7 Gradient of the ELBO Loss:

$$\nabla_\lambda L(\lambda) = \nabla_\lambda \text{ELBO}(\lambda) \quad (9.8)$$

$$= \nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y, \theta)] + H(q_\lambda) \quad (9.8)$$

$$= \nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] - \text{KL}(q_\lambda(\theta) \parallel p(\theta)) \quad (9.8)$$

$$= \nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] - \nabla_\lambda \text{KL}(q_\lambda(\theta) \parallel p(\theta)) \quad (9.8)$$

Problem

In order to use SGD we need to evaluate the gradient of the loss:

$$\nabla_\lambda \mathbb{E}_{\theta \sim p} [l(\theta; \mathbf{x})] = \mathbb{E} [\nabla_{\mathbf{x} \sim p} l(\theta; \mathbf{x})] = \frac{1}{m} \sum_{i=1}^m \nabla_{\mathbf{x}_i \sim p} l(\theta; \mathbf{x}_i)$$

however in eq. (9.8) only the second term can be derived easily. For the first term we cannot move the gradient inside the expectation as the expectations depends on the parameter w.r.t. which we differentiate:

$$\nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] = \frac{\partial}{\partial \lambda} \int q_\lambda \log p(y|\theta) d\theta$$

Solutions:

- Score Gradients
- Reparameterization Trick: reparameterize a function s.t. it depends on another parameter and reformulate it s.t. it still returns the same value.

4. The Reparameterization Trick

Principle 9.1

[proof 10.14]

Reparameterization Trick: Let ϕ be some base distribution from which we can sample and assume there exist an invertible function g s.t. $\theta = g(\epsilon, \lambda)$ then we can write θ in terms of a new distribution parameterized by $\epsilon \sim \phi(\epsilon)$:

$$\theta \sim q(\theta|\lambda) = \phi(\epsilon)|\nabla_\epsilon g(\epsilon; \lambda)|^{-1} \quad (9.9)$$

we can then write by the law of the unconscious statistician law 65.6:

$$\mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] = \mathbb{E}_{\epsilon \sim \phi} [\log p(y|g(\epsilon; \lambda))] \quad (9.10)$$

\Rightarrow the expectations does not longer depend on λ and we can pull in the gradient!

$$\nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] = \nabla_\lambda \mathbb{E}_{\epsilon \sim \phi} [\log p(y|g(\epsilon; \lambda))] \quad (9.11)$$

$$= \mathbb{E}_{\epsilon \sim \phi} [\nabla_\lambda \log p(y|g(\epsilon; \lambda))] \quad (9.12)$$

Definition 9.8

[example 10.3]

Reparameterized ELBO Gradient^[def. 9.7]:

By using the reparameterization trick principle 9.1 we can write the gradient of the ELBO as:

$$\begin{aligned} \nabla_\lambda L(\lambda) &= \nabla_\lambda \text{ELBO}(\lambda) \\ &= \nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] - \nabla_\lambda \text{KL}(q_\lambda(\theta) \parallel p(\theta)) \\ &= \mathbb{E}_{\epsilon \sim \phi} [\nabla_\lambda \log p(y|g(\epsilon; \lambda))] - \nabla_\lambda \text{KL}(q_\lambda(\theta) \parallel p(\theta)) \end{aligned} \quad (9.13)$$

Corollary 9.2

[proof 10.3]

Reparameterized ELBO for Gaussians:

Lets assume a Gaussian distribution for our approximate distribution: q and lets use a normal distribution for $\phi(\epsilon)$:
 $\theta \sim q(\theta|\lambda) = \mathcal{N}(\theta; \mu, \Sigma) \Rightarrow \lambda = [\mu \Sigma]$
 $\epsilon \sim \phi(\epsilon) = \mathcal{N}(\epsilon; 0, I)$

Then it follows that the ELBO:

$$\begin{aligned} \nabla_\lambda L(\lambda) &= \nabla_\lambda \text{ELBO}(\lambda) \\ &= \nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] - \nabla_\lambda \text{KL}(q_\lambda(\theta) \parallel p(\theta)) \\ &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\nabla_{\mathbf{C}, \mu} \log p(y|g(\epsilon; \mathbf{C}, \mu))] \\ &\quad - \nabla_{\mathbf{C}, \mu} \text{KL}(q_\lambda(\mathbf{C}, \mu) \parallel p(\theta)) \\ &\approx \frac{n}{m} \sum_{j=1}^m \nabla_{\mathbf{C}, \mu} \log p(y_j | g(\epsilon_j; \mathbf{C}, \mu) + \mu, \mathbf{x}_j) \\ &\quad - \nabla_{\mathbf{C}, \mu} \text{KL}(q_\lambda(\mathbf{C}, \mu) \parallel p(\theta)) \end{aligned} \quad (9.14)$$

2. Markov Chain Monte Carlo Methods

Definition 9.9

Markov Chain Monte Carlo (MCMC) Methods:

3. Integrated Nested Laplace Approximation

$$\eta_i = \alpha + \sum_{j=1}^{n_f} f^{(j)}(\mathbf{u}_{ji}) + \sum_{k=1}^{n_\beta} \beta_k z_{ki} + \epsilon_i \quad (9.15)$$

$$p(\mathbf{x}, \theta)p(\mathbf{y}) = p(\mathbf{x}) \quad (9.16)$$

$$p(\mathbf{x}_i|\mathbf{y}) = \int p(\mathbf{x}_i|\theta, \mathbf{y})p(\theta|\mathbf{y}) d\theta$$

$$\rightarrow \tilde{p}(\mathbf{x}_i|\mathbf{y}) = \int \tilde{p}(\mathbf{x}_i|\theta, \mathbf{y}) \tilde{p}(\theta|\mathbf{y}) d\theta$$

$$p(\theta_j|\mathbf{y}) = \int p(\theta|\mathbf{y}) d\theta_{-j}$$

$$\rightarrow \tilde{p}(\theta_j|\mathbf{y}) = \int \tilde{p}(\theta|\mathbf{y}) d\theta_{-j}$$

$p(\mathbf{x}_i|\theta, \mathbf{y})$ and $p(\theta|\mathbf{y})$ are approximated and the posterior marginal densities are then calculated using numerical integration:

Note

The numerical integration is possible if θ is small i.e. $m = \dim(\theta) \leq 5$.

4. Approximating $p(\theta|\mathbf{y})$ and $p(\mathbf{x}_i|\mathbf{y})$

$$p(\mathbf{x}, \theta, \mathbf{y}) = p(\mathbf{x}|\theta, \mathbf{y})p(\theta, \mathbf{y}) = p(\mathbf{x}|\theta, \mathbf{y})\tilde{p}(\theta|\mathbf{y})p(\mathbf{y})$$

$$\Rightarrow \tilde{p}(\theta|\mathbf{y}) = \frac{p(\mathbf{x}, \theta, \mathbf{y})}{\tilde{p}(\mathbf{x}|\theta, \mathbf{y})p(\mathbf{y})} \propto \left. \frac{p(\mathbf{x}, \theta, \mathbf{y})}{p_G(\mathbf{x}|\theta, \mathbf{y})} \right|_{\mathbf{x}=\mathbf{x}^*(\theta)}$$

1. Marginal Posterior of the latent field $p(\mathbf{x}_i|\mathbf{y})$ are calculated by first approximating $p(\theta|\mathbf{y})$:

$$p(\theta|\mathbf{y})_G = \mathcal{N}(x_i; \mu_i(\theta), \sigma_i^2(\theta))$$

and then numerical integration w.r.t. θ :

$$\tilde{p}(\mathbf{x}_i|\mathbf{y}) = \sum_k P_G(\theta_k|\mathbf{y}) \tilde{p}(\theta_k|\mathbf{y}) \Delta_k$$

Note

$\tilde{p}(\theta|\mathbf{y})$ is usually quite different from a Gaussian s.t. the Gaussian approximation alone is not really sufficient.

Bayesian Neural Networks (BNN)

Definition 10.1 Bayesian Neural Networks (BNN):
 ① Model the prior over our weights $\theta = [W^0 \dots W^L]$ by a neural network:

$$\theta \sim p_{\lambda}(\theta) = F \quad \text{with} \quad F^L = \varphi \circ \bar{F}^L = \varphi(W^L x + b^L)$$

for each weight $w^{(0)}_{k,j}$ of input x_j with weight on the hidden variable $z_i^{(0)}$ with $a_i^0 = \varphi(z_i^{(0)})$ it follows:

$$w^{(0)}_{k,j} = p_w(a_{k,j}) \text{ i.e. } \mathcal{N}(\mu_{k,j}, \sigma_{k,j}^2)$$

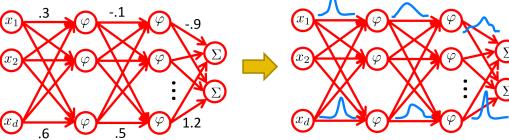


Figure 8

② The parameters of likelihood function are modeled by the output of the network:

$$p(y|F(\theta, X)) \quad \text{see example 10.4} \quad (10.1)$$

Note

Recall for normal Bayesian Linear regression we had:

Problem

All the weights of the prior $p_{\lambda}(\theta) = F$ are correlated in some complex way see Figure 8. Thus even if the prior and likelihood are simple, the posterior will be not. \Rightarrow need to approximate the posterior $p(\theta|y, X)$ i.e. by fitting a Gaussian distribution to each weight of the posterior neural network.

0.0.1. MAP estimates for BNN

Definition 10.2 BNN MAP Estimate:

We need to do a forward pass for each x_i in order to obtain $\mu(x_i; \theta)$ and $\sigma(x_i; \theta)^2$:

$$\theta^* = \arg \max_{\theta} \{p(\theta|X, y)\} \stackrel{\text{eq. (6.13)}}{=} \arg \min_{\theta} \|\lambda\|_2^2 \\ - \sum_{i=1}^n \left(\frac{1}{2\sigma(x_i; \theta)^2} \|y_i - \mu(x_i; \theta)\|^2 + \frac{1}{2} \log \sigma(x_i; \theta)^2 \right)$$

Explanation 10.1. [def. 10.2]

- $\frac{1}{2} \log \sigma(x_i; \theta)^2$: tries to force neural network to predict small uncertainty
- $\frac{1}{2\sigma(x_i; \theta)^2} \|y_i - \mu(x_i; \theta)\|^2$: tries to force neural network to predict accurately but if this is not possible for certain data points the network can attenuate the loss to a larger variance.

Definition 10.3 proof 10.15

MAP Gradient of BNN:

$$\theta_{t+1} = \theta_t (1 - 2\lambda \eta_t) - \eta_t \nabla \sum_{i=1}^n \log p(y_i|x_i, \theta) \quad (10.2)$$

Note

- The gradients of the objective eq. (10.2) can be calculated using auto-differentiation techniques e.g. Pytorch or Tensorflow.
- The BNN MAP estimate fails to predict epistemic uncertainty^[def. 6.16] \iff it is overconfident in regions where we haven't even seen any data.
 \Rightarrow need to use Bayesian approach to approximate posterior distribution.

1. Variational Inference For BNN

We use the objective eq. (9.14) as loss in order to perform back propagation.

2. Making Predictions

Proposition 10.1 Title:

1. Proofs

$$\begin{aligned} \text{Proof 10.1: Definition 6.15:} \\ p(f_*|x_*, X, y) &= \frac{p(f_*, x_*, X, y)}{p(x_*, X, y)} \\ &= \frac{\int p(f_*, x_*, X, y, w) dw}{p(x_*, X, y)} \\ \stackrel{\text{eq. (65.19)}}{=} &\frac{\int p(f_*|x_*, X, y, w) p(w|x_*, X, y) dw}{p(x_*, X, y)} \\ \stackrel{\text{eq. (65.19)}}{=} &\int p(f_*|x_*, X, y, w) p(w|X, y) dw \\ &\stackrel{+}{=} \int p(f_*|x_*, w) p(w|X, y) dw \end{aligned}$$

Note ♣

- f_* is independent of $\mathcal{D} = \{X, y\}$ given the fixed parameter w .
- w does only depend on the observed data $\mathcal{D} = \{X, y\}$ and not the unseen data x_* .

Proof 10.2: Definition 6.13:

$$\begin{aligned} p(y|X) &= \int p(y, w|X) dw = \int p(y|w, X)p(w|X) dw \\ \stackrel{\text{eq. (6.6)}}{=} &\int p(y|w, X)p(w) dw \end{aligned}$$

Proof 10.3: Definition 7.4:

$$\begin{aligned} p(x_t, x_{t-1}|y_{1:t}) &\stackrel{\text{eq. (65.19)}}{=} p(x_t|x_{t-1}, y_{1:t})p(x_{t-1}|y_{1:t}) \\ &\stackrel{\text{independ.}}{=} p(x_t|x_{t-1})p(x_{t-1}|y_{1:t}) \end{aligned}$$

marginalization/integration over x_{t-1} gives the desired result.

Proof 10.4: Definition 7.5:

$$\begin{aligned} p(x_t, y_t|y_{1:t-1}) &\stackrel{\text{eq. (65.23)}}{=} \left\{ \begin{array}{l} p(x_t|y_t, y_{1:t-1})p(y_t|y_{1:t-1}) \\ p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) \end{array} \right. \\ &\dots \\ p(y_t|x_t, y_{1:t-1}) &\stackrel{[\text{cor. 7.2}]}{=} p(y_t|x_t) \\ &\dots \end{aligned}$$

from which follows immediately eq. (7.5).

Proof 10.5: Definition 7.6:

$$\begin{aligned} p(y_t|y_{1:t-1}) &= \int p(y_t, x_t|y_{1:t-1}) dx_t \\ &= \int p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1}) dx_t \\ \stackrel{[\text{cor. 7.2}]}{=} &\int p(y_t|x_t)p(x_t|y_{1:t-1}) dx_t \end{aligned}$$

Proof 10.6: [cor. 7.3]:

$$\begin{aligned} p(x_{1:t}, y_{1:t}) &\stackrel{\text{eq. (65.19)}}{=} p(y_{1:t}|x_{1:t})p(x_{1:t}) \\ &\stackrel{\text{law 65.2}}{=} p(y_{1:t}|x_{1:t})p(x_t|x_{t-1:0}) \cdots p(x_2|x_1)p(x_1) \\ &\stackrel{\text{eq. (7.1)}}{=} p(y_{1:t}|x_{1:t}) \left(p(x_1) \prod_{i=2}^t p(x_i|x_{i-1}) \right) \\ &\stackrel{\text{law 65.2}}{\stackrel{[\text{cor. 7.2}]}{=}} \left(p(y_1|x_1) \cdots p(y_t|x_t) \right) \left(p(x_1) \prod_{i=2}^t p(x_i|x_{i-1}) \right) \\ &= \underbrace{p(y_1|x_1)p(x_1)}_{= 1} \prod_{i=2}^t p(y_i|x_i)p(x_i|x_{i-1}) \end{aligned}$$

Proof 10.7: GP Posterior Distribution^[def. 8.1]

$$\begin{aligned} p(w|\mathcal{D}) &\propto p(\mathcal{D}|w)p(w) \\ &\propto \exp \left(-\frac{1}{2} \frac{1}{\sigma_n^2} (y - Xw)^T (y - Xw) \right) \exp \left(-\frac{1}{2} w^T \Sigma^{-1} w \right) \\ &\propto \exp \left\{ -\frac{1}{2} \frac{1}{\sigma_n^2} (y^T y - 2w^T X^T y + w^T X^T X w + \sigma_n^2 w^T \Sigma^{-1} w) \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \frac{1}{\sigma_n^2} (y^T y - 2w^T X^T y + w^T (X^T X + \sigma_n^2 \Sigma^{-1}) w) \right\} \end{aligned}$$

We know that a Gaussian $\mathcal{N}(w|\mu_w, \Sigma_w^{-1})$ should look like:

$$\begin{aligned} p(w|\mathcal{D}) &\propto \exp \left(-\frac{1}{2} (w - \mu_w)^T \Sigma_w (w - \mu_w) \right) \\ &\propto \exp \left(-\frac{1}{2} \left(w^T \Sigma_w w - 2w^T \Sigma_w \mu_w + \mu_w^T \Sigma_w \mu_w \right) \right) \\ \Sigma_w &\text{ follows directly } \Sigma_w = \sigma_n^{-2} X^T X + \Sigma_p \\ \mu_w &\text{ follows from } 2w^T X^T y = 2w^T \Sigma_w \mu_w \Rightarrow \mu_w = \Sigma_w^{-1} X^T y. \end{aligned}$$

Proof 10.8: [def. 8.2]

Proof 10.9: [def. 9.4] In a Bayesian setting we are usually interested in maximizing the log prior+likelihood:

$$\begin{aligned} L_n(\theta) &= \log(p(\theta|y)) = (\text{log Prior} + \text{log Likelihood}) \\ \text{we now approximate } L_n(\theta) &\text{ by a Taylor approximation around its maximum } \hat{\theta}: \\ L_n(\theta) &= L_n(\hat{\theta}) + \frac{1}{2} \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} (\theta - \hat{\theta})^2 + \mathcal{O}((\theta - \hat{\theta})^3) \end{aligned}$$

we can no derive the distribution:

$$\begin{aligned} p(\theta|y) &\approx \exp(L_n(\theta)) = \exp(\log p(\theta|y)) \\ &= p(\hat{\theta}) \exp \left(\frac{1}{2} \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} \right) \\ &= \sqrt{2\pi\sigma^2} p(\hat{\theta}) \mathcal{N}(\theta; \hat{\theta}, \sigma) \approx \frac{1}{\sqrt{2\pi\sigma^2}} \mathcal{N}(\theta; \hat{\theta}, \sigma) \end{aligned}$$

Notes

- the derivative of the maximum must be zero by definition $\frac{\partial L_n}{\partial \theta} \Big|_{\hat{\theta}} = 0$
- we approximate the normalization constant $\frac{1}{Z}$ by $\sqrt{2\pi\sigma^2} p(\hat{\theta})$.

Proof 10.10: [def. 9.4] 2D:

$$\begin{aligned} \nabla L_n(\theta) &= \nabla L_n(\theta_1, \theta_2) = 0 \\ L_n(\theta) &= L_n(\hat{\theta}) + \frac{1}{2} (A(\theta_1 - \hat{\theta}_1)^2 + B(\theta_2 - \hat{\theta}_2)^2 \\ &\quad + C(\theta_1 - \hat{\theta}_1)(\theta_2 - \hat{\theta}_2)) \end{aligned}$$

$$\begin{aligned} L_n(\theta) &= L_n(\hat{\theta}) + (\theta - \hat{\theta})^T H(\hat{\theta})(\theta - \hat{\theta}) \\ &= L_n(\hat{\theta}) + \frac{1}{2} Q(\theta) \end{aligned}$$

$$\begin{aligned} A &= \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} & B &= \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} & C &= \frac{\partial^2 L_n}{\partial \theta_1 \partial \theta_2} \Big|_{\hat{\theta}} \\ H &= \begin{bmatrix} A & C \\ C & B \end{bmatrix} & \Sigma &= H^{-1}(\hat{\theta}) \end{aligned}$$

Proof 10.11: [def. 9.4] k-dimensional:

$$\begin{aligned} L_n(\theta) &\approx L_n(\hat{\theta}) + (\theta - \hat{\theta})^T \nabla \nabla^T L_n(\hat{\theta})(\theta - \hat{\theta}) \\ H(\theta) &= \nabla \nabla^T L_n(\theta) & \Sigma &= H^{-1}(\hat{\theta}) \\ p(\theta|y) &= \sqrt{(2\pi)^n \det(\Sigma)} p(\hat{\theta}) \mathcal{N}(\theta; \hat{\theta}, \Sigma) \\ &\approx c \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \mathcal{N}(\theta; \hat{\theta}, \Sigma) \end{aligned}$$

Proof 10.12: [def. 9.5]

$$\begin{aligned}
 q^* &\in \arg \min_{q \in Q} \text{KL}(q(\theta) \| p(\theta|y)) \\
 p(\theta|y) &= \frac{1}{Z} p(\theta, y) \\
 &= \arg \min_q \mathbb{E}_{\theta \sim q} \left[\log \frac{q(\theta)}{\frac{1}{Z} p(\theta, y)} \right] \\
 &= \arg \min_q \mathbb{E}_{\theta \sim q} \left[\log q(\theta) - \log \frac{1}{Z} - \log p(\theta, y) \right] \\
 &= \arg \min_q \mathbb{E}_{\theta \sim q} \left[-\log q(\theta) + \mathbb{E}_{\theta \sim q} [\log Z] \right] \\
 &\quad - \mathbb{E}_{\theta \sim q} [\log p(\theta, y)] \\
 &= \arg \max_q \mathbb{E}_{\theta \sim q} [\log p(\theta, y)] + H(q) \\
 &= \arg \max_q \mathbb{E}_{\theta \sim q} [\log p(\theta|y) + \log p(\theta) - \log q(\theta)] \\
 &= \arg \max_q \mathbb{E}_{\theta \sim q} [\log p(\theta|y) + \text{KL}(q(\theta) \| p(\theta))]
 \end{aligned}$$

Proof 10.13: [def. 9.6]

$$\begin{aligned}
 \log p(y) &= \log \int p(y, \theta) d\theta = \log \int p(y|\theta)p(\theta) d\theta \\
 &= \log \int p(y|\theta) \frac{p(\theta)}{q_\lambda(\theta)} q_\lambda(\theta) d\theta \\
 &= \log \mathbb{E}_{\theta \sim q_\lambda} \left[p(y|\theta) \frac{p(\theta)}{q_\lambda(\theta)} \right] \\
 \text{eq. (65.55)} &\geq \mathbb{E}_{\theta \sim q_\lambda} \left[\log \left(p(y|\theta) \frac{p(\theta)}{q_\lambda(\theta)} \right) \right] \\
 &= \mathbb{E}_{\theta \sim q_\lambda} \left[\log p(y|\theta) - \log \frac{p(\theta)}{q_\lambda(\theta)} \right] \\
 &= \mathbb{E}_{\theta \sim q_\lambda} [\log p(y|\theta)] - \text{KL}(q_\lambda \| p(\cdot))
 \end{aligned}$$

Proof 10.14: principle 9.1 Let:

$$\begin{aligned}
 \epsilon &\sim \phi(\epsilon) \quad \text{correspond to} \quad X \sim f_X \\
 \theta &= g(\epsilon; \lambda) \quad \mathcal{Y} = \{y | y = g(x), \forall x \in \mathcal{X}\} \\
 \text{then it follows immediately with formula 65.2:} \\
 \theta &\sim q_\lambda(\theta) = q(\theta|\lambda) = \frac{f_X(g^{-1}(y))}{\left| \frac{dg}{dx}(g^{-1}(y)) \right|} \\
 &= \phi(\epsilon) |\nabla_\epsilon g(\epsilon; \lambda)|^{-1}
 \end{aligned}$$

\Rightarrow parameterized in terms of ϵ

Proof 10.15: [def. 10.3]

$$\begin{aligned}
 \theta_{t+1} &= \theta_t - \eta_t \left(\nabla \log p(\theta) - \nabla \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \theta) \right) \\
 &= \theta_t - \eta_t \left(2\lambda \theta_t - \nabla \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \theta) \right) \\
 &= \theta_t (1 - 2\lambda \eta_t) - \eta_t \nabla \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \theta)
 \end{aligned}$$

Example 10.2 ELBO Bayesian Logistic Regression:

Suppose:

$$\begin{aligned}
 Q &= \text{diag. Gaussians} \quad \Rightarrow \quad \lambda = [\mu_{1:d} \ \ \sigma_{1:d}^2] \in \mathbb{R}^{2d} \\
 p(\theta) &= \mathcal{N}(0, I)
 \end{aligned}$$

Then it follows for the terms of the ELBO:

$$\begin{aligned}
 \text{KL}(q_\lambda \| p(\theta)) &= \frac{1}{2} \sum_{i=1}^d (\mu_i^2 + \sigma_i^2 - 1 - \ln \sigma_i^2) \\
 \mathbb{E}_{\theta \sim q_\lambda} [p(y|\theta)] &= \mathbb{E}_{\theta \sim q_\lambda} \left[\sum_{i=1}^n \log p(y_i | \theta, \mathbf{x}_i) \right] \\
 &= \mathbb{E}_{\theta \sim q_\lambda} \left[-\sum_{i=1}^n \log (1 + \exp(-y_i \theta^\top \mathbf{x}_i)) \right]
 \end{aligned}$$

Example 10.3 ELBO Gradient Gaussian: Suppose:

$$\begin{aligned}
 \theta &\sim q(\theta|\lambda) = \mathcal{N}(\theta; \mu, \Sigma) \quad \Rightarrow \quad \lambda = [\mu \ \ \Sigma] \\
 \epsilon &\sim \phi(\epsilon) = \mathcal{N}(\epsilon; 0, I)
 \end{aligned}$$

we can reparameterize using principle 9.1 by using:

$$\theta \sim g(\epsilon; \lambda) = \mathbf{C}\epsilon + \mu \quad \text{with} \quad \mathbf{C} : \mathbf{CC}^\top = \Sigma$$

from this it follows: (\mathbf{C} is the Cholesky factor of Σ)

$$g^{-1}(\theta, \lambda) = \epsilon = \mathbf{C}^{-1}(\theta - \mu) \quad \frac{\partial g(\epsilon; \lambda)}{\partial \epsilon} = \mathbf{C}$$

from this it follows:

$$\begin{aligned}
 q(\theta|\lambda) &= \frac{\phi(\epsilon)}{\left| \frac{dg(\epsilon; \theta)}{d\epsilon} (g^{-1}(\theta)) \right|} = \phi(\epsilon) |C|^{-1} \\
 \Leftrightarrow \phi(\epsilon) &= q(\theta|\lambda) |C|
 \end{aligned}$$

we can then write the reparameterized expectation part of the gradient of the ELBO as:

$$\begin{aligned}
 \nabla_\lambda L(\lambda)_1 &= \nabla_\lambda \mathbb{E}_{\epsilon \sim \phi} [\log p(y|g(\epsilon; \lambda))] \\
 &= \nabla_{\mathbf{C}, \mu} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\log p(y|\mathbf{C}\epsilon + \mu)] \\
 \text{i.i.d.} &\quad \nabla_{\mathbf{C}, \mu} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[\sum_{i=1}^n \log p(y_i | \mathbf{C}\epsilon + \mu, \mathbf{x}_i) \right] \\
 &= \nabla_{\mathbf{C}, \mu} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[n \frac{1}{n} \sum_{i=1}^n \log p(y_i | \mathbf{C}\epsilon + \mu, \mathbf{x}_i) \right] \\
 &= \nabla_{\mathbf{C}, \mu} n \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\mathbb{E}_{i \sim \mathcal{U}(\{1, n\})} \log p(y_i | \mathbf{C}\epsilon + \mu, \mathbf{x}_i)]
 \end{aligned}$$

Draw a mini batch $\begin{cases} \epsilon^{(1)}, \dots, \epsilon^{(m)} \\ j_1, \dots, j_m \sim \mathcal{U}(\{1, n\}) \end{cases}$

$$= n \frac{1}{m} \sum_{j=1}^m \nabla_{\mathbf{C}, \mu} \log p(y_j | \mathbf{C}\epsilon + \mu, \mathbf{x}_j)$$

$$\nabla_\lambda L(\lambda) = \nabla_\lambda \text{ELBO}(\lambda) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\nabla_{\mathbf{C}, \mu} \log p(y|\mathbf{C}\epsilon + \mu)] - \nabla_{\mathbf{C}, \mu} (q_{\mathbf{C}, \mu} \| p(\theta))$$

Example 10.4 BNN Likelihood Function Examples:

$$p(y|\mathbf{X}, \theta) = \begin{cases} \mathcal{N}(y; \mathbf{F}(\mathbf{X}, \theta), \sigma^2) \\ \mathcal{N}(y; \mathbf{F}(\mathbf{X}, \theta)_1, \exp \mathbf{F}(\mathbf{X}, \theta)_1) \end{cases}$$

2. Examples

Example 10.1 Laplace Approximation

Logistic Regression Likelihood + Gaussian Prior:

Kernels

Given objects we cannot assume that they are vectors/can be represented as vectors in feature space.

Hence it is also not guaranteed that those objects can be added and multiplied by scalars.

Question: then how can we define a more general notion of similarity?

Definition 11.1 Similarity Measure $\text{sim}(A, B)$: A similarity measure or similarity function is a real-valued function that quantifies the similarity between two objects.

No single definition of a similarity measure exists but often they are defined in terms of the inverse of distance metrics and they take on large values for similar objects and either zero or a negative value for very dissimilar objects.

Definition 11.2 Dissimilarity Measure $\text{dissim}(A, B)$: Is a measure of how dissimilar objects are, rather than how similar they are.

Thus it takes the largest values for objects that are really far apart from another.

Dissimilarities are often chosen as the squared norm of two difference vectors:

$$\begin{aligned} \|\mathbf{x} - \mathbf{y}\|^2 &= \mathbf{x}^\top \mathbf{x} + \mathbf{y}^\top \mathbf{y} - 2\mathbf{x}^\top \mathbf{y} \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d \\ \text{dissim}(\mathbf{x}, \mathbf{y}) &= \text{sim}(\mathbf{x}, \mathbf{x}) + \text{sim}(\mathbf{y}, \mathbf{y}) - 2\text{dissim}(\mathbf{x}, \mathbf{y}) \end{aligned} \quad (11.1)$$

Attention

It is better to rely on similarity measures instead of dissimilarity measures. Dissimilarities are often not adequate from a modeling point of view, because for objects that are really dissimilar/far from each other, we usually have the biggest problem to estimate their distance.

E.g. for a bag of words it is easy to determine similar words, but it is hard to estimate which words are most dissimilar. For normed vectors the only information of a dissimilarity defined as in eq. (11.1) becomes $2\mathbf{x}^\top \mathbf{y} = 2\text{dissim}(\mathbf{x}, \mathbf{y})$

Definition 11.3 Feature Map ϕ : is a mapping $\phi : \mathcal{X} \mapsto \mathcal{V}$ that takes an input $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$ and maps it into another feature space $\mathcal{V} \subseteq \mathbb{R}^D$.

Note

Such feature maps can lead to an exponential number of terms i.e. for a polynomial feature map, with monomials of degree up to p and feature vectors of dimension $\mathbf{x} \in \mathbb{R}^d$ we obtain a feature space of size:

$$D = \dim(\mathcal{V}) = \binom{p+d}{d} = \mathcal{O}(d^p) \quad (11.2)$$

when using the polynomial kernel^[def. 11.10], this can be reduced to the order d .

Definition 11.4 Kernel \mathbf{k} : Let $\mathcal{X} \subseteq \mathbb{R}^d$ be the data space. A map $\mathbf{k} : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is called kernel if there exists an inner product space^[def. 59,78] called **feature space** $(\mathcal{V}, \langle \cdot, \cdot \rangle_{\mathcal{V}})$ and a map $\phi : \mathcal{X} \mapsto \mathcal{V}$ s.t.

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{V}} \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X} \quad (11.3)$$

Corollary 11.1 Kernels and similarity: Kernels are defined in terms of inner product spaces and hence have a notion of similarity between its arguments.

Example

Let $\mathbf{k}(\mathbf{x}, \mathbf{y}) := \mathbf{x}^\top \mathbf{A} \mathbf{y}$ thus the kernel measures the similarity between \mathbf{x} and \mathbf{y} by the inner product $\mathbf{x}^\top \mathbf{y}$ weighted by the matrix \mathbf{A} .

Corollary 11.2 Kernels and distance: Let $\mathbf{k}(\mathbf{x}, \mathbf{y})$ be a measure of similarity between \mathbf{x} and \mathbf{y} then \mathbf{k} induces a dissimilarity/distance between \mathbf{x} and \mathbf{y} defined as the difference between the self-similarities $\mathbf{k}(\mathbf{x}, \mathbf{x}) + \mathbf{k}(\mathbf{y}, \mathbf{y})$ and the cross-similarities $\mathbf{k}(\mathbf{x}, \mathbf{y})$:

$$\text{dissimilarity}(\mathbf{x}, \mathbf{y}) := \mathbf{k}(\mathbf{x}, \mathbf{x}) + \mathbf{k}(\mathbf{y}, \mathbf{y}) - 2\mathbf{k}(\mathbf{x}, \mathbf{y})$$

Note

The factor 2 is required to ensure that $d(\mathbf{x}, \mathbf{x}) = 0$.

1. The Gram Matrix

Definition 11.5 Kernel (Gram) Matrix:

Given: a mapping $\phi : \mathbb{R}^d \mapsto \mathbb{R}^D$ and a corresponding kernel function $\mathbf{k} : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ with $\mathbf{k} \subseteq \mathbb{R}^d$.

Let S be any finite subset of data $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathcal{X}$. Then the kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ is defined by:

$$\begin{aligned} \mathbf{K} &= \phi(\mathbf{X})\phi(\mathbf{X})^\top = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))(\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))^\top \\ &= \begin{pmatrix} \mathbf{k}(\mathbf{x}_1, \mathbf{x}_1) & \dots & \mathbf{k}(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ \mathbf{k}(\mathbf{x}_n, \mathbf{x}_1) & \dots & \mathbf{k}(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix} = \begin{pmatrix} \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_n) \end{pmatrix} \\ \mathbf{K}_{ij} &= \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) \end{aligned}$$

Corollary 11.3

VΛV^T

Kernel Eigenvector Decomposition:

For any symmetric matrix (Gram matrix $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)_{i,j=1}^n$) there exists an eigenvector decomposition:

$$\mathbf{K} = \mathbf{V}\Lambda\mathbf{V}^\top \quad (11.4)$$

\mathbf{V} : orthogonal matrix of eigenvectors $(\mathbf{v}_t, i)_{i=1}^n$

Λ : diagonal matrix of eigenvalues λ_i

Assuming all eigenvalues λ_t are non-negative, we can calculate the mapping:

$$\phi : \mathbf{x}_i \mapsto (\sqrt{\lambda_t} \mathbf{v}_t, i)_{t=1}^n \in \mathbb{R}^n, \quad i = 1, \dots, n \quad (11.5)$$

which allows us to define the Kernel \mathbf{K} as:

$$\begin{aligned} \phi^\top(\mathbf{x}_i)\phi(\mathbf{x}_j) &= \sum_{t=1}^n \lambda_t \mathbf{v}_{t,i} \mathbf{v}_{t,j} = (\mathbf{V}\Lambda\mathbf{V}^\top)_{i,j} = \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) \end{aligned} \quad (11.6)$$

1. Necessary Properties

Property 11.1 Inner Product Space:

\mathbf{k} must be an inner product of a suitable space \mathcal{V} .

Property 11.2 Symmetry:

$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \mathbf{k}(\mathbf{y}, \mathbf{x}) = \phi(\mathbf{x})^\top \phi(\mathbf{y}) = \phi(\mathbf{y})^\top \phi(\mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$

Property 11.3 Non-negative Eigenvalues/p.s.d.s Form:

Let $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ be an n -set of a finite input space \mathcal{V} . A kernel \mathbf{k} must induces a p.s.d. symmetric kernel matrix \mathbf{K} for any possible $S \subseteq \mathcal{X}$ see ?? 11.11.

all eigenvalues of the kernel gram matrix \mathbf{K} for finite \mathcal{V} must be non-negative ?? 59.2.

Notes

- The extension to infinite dimensional Hilbert Spaces might also include a non-negative weighting/eigenvalues:

$$\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{z})$$

- In order to be able to use a kernel, we need to verify that the kernel is p.s.d. for all n-vectors $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, as well as for future unseen values.

2. Mercers Theorem

Theorem 11.1 Mercers Theorem:

Let \mathcal{X} be a compact subset of \mathbb{R}^n and $\mathbf{k} : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ a kernel function.

Then one can expand \mathbf{k} in a uniformly convergent series of bounded functions ϕ s.t.

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{x}') \quad (11.7)$$

Theorem 11.2 General Mercers Theorem:

Let Ω be a compact subset of \mathbb{R}^n . Suppose \mathbf{k} is a gernal continuous symmetric function such that the integral operator:

$$T_{\mathbf{k}} : L_2(\mathbf{X}) \mapsto L_2(\mathbf{X}) \quad (T_{\mathbf{k}} f)(\cdot) = \int_{\Omega} \mathbf{k}(\cdot, \mathbf{x}) f(\mathbf{x}) d\mathbf{x} \quad (11.8)$$

is positive, that is it satisfies:

$$\int_{\Omega \times \Omega} \mathbf{k}(\mathbf{x}, \mathbf{z}) f(\mathbf{x}) f(\mathbf{z}) d\mathbf{x} d\mathbf{z} > 0 \quad \forall f \in L_2(\Omega)$$

Then we can expand $\mathbf{k}(\mathbf{x}, \mathbf{z})$ in a uniformly convergent series in terms of $T_{\mathbf{k}}$'s eigen-functions $\phi_j \in L_2(\Omega)$, with $\|\phi_j\|_{L_2} = 1$ and positive associated eigenvalues $\lambda_j > 0$.

Note

All kernels satisfying mercers conditions describe an inner product in a high dimensional space.

\Rightarrow can replace the inner product by the kernel function.

3. The Kernel Trick

Definition 11.6 Kernel Trick:

If a function has an analytic form we do no longer need to calculate:

- the function mapping $\mathbf{x} \mapsto \phi(\mathbf{x})$ and

explicitly but simply us the formula for the kernel:

$$\phi(\mathbf{x})^\top \phi(\mathbf{y}) = \mathbf{k}(\mathbf{x}, \mathbf{y}) \quad (11.9)$$

see examples 11.1 and 11.2

Note

If we chose h small, all data points not close to h will be 0/discarded \iff data points are considered as independent. Length of all vectors in feature space is one $\mathbf{k}(\mathbf{x}, \mathbf{x}) = e^0 = 1$. Thus: Data points in input space are projected onto a high-(infinity-)dimensional sphere in feature space.

Classification: Cutting with hyperplanes through the sphere. How to chose h : good heuristics, take median of the distance all points but better is cross validation.

6. The Matern Kernel

When looking at actual data/sample paths the smoothness of the Gaussian kernel^[def. 11.13] is often a too strong assumption that does not model reality the same holds true for the non-smoothness of the exponential kernel^[def. 11.12]. A solution to this dilemma is the Matern kernel.

Definition 11.14 Matern Kernel: is a kernel which allows you to specify the level of smoothness $\mathbf{k} \in \mathcal{C}^{[\nu]}$ by a positive parameter ν :

$$\mathbf{k}(x, y) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\rho} \right)^{\nu} \mathbf{K}_{\nu} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\rho} \right) \quad (11.17)$$

$\nu, \rho \in \mathbb{R}_+$ ν : Smoothness
 \mathbf{K}_{ν} modified Bessel function of the second kind

6. Kernel Engineering

Often linear and even non-linear simple kernels are not sufficient to solve certain problems, especially for pairwise problems i.e. user & product, exon & intron,... Composite kernels can be the solution to such problems.

1. Closure Properties/Composite Rules

Suppose we have two kernels:

$$\mathbf{k}_1 : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R} \quad \mathbf{k}_2 : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$$

defined on the data space $\mathcal{X} \subseteq \mathbb{R}^d$. Then we may define using Composite Rules:

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \mathbf{k}_1(\mathbf{x}, \mathbf{x}') + \mathbf{k}_2(\mathbf{x}, \mathbf{x}') \quad (11.18)$$

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \mathbf{k}_1(\mathbf{x}, \mathbf{x}') \cdot \mathbf{k}_2(\mathbf{x}, \mathbf{x}') \quad (11.19)$$

$$\alpha \in \mathbb{R}_+$$

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \alpha \mathbf{k}_1(\mathbf{x}, \mathbf{x}') \quad (11.20)$$

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = f(\mathbf{x}) f(\mathbf{x}') \quad (11.21)$$

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \mathbf{k}_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (11.22)$$

$$\alpha \in \mathbb{R}_+$$

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = p(\mathbf{k}(\mathbf{x}, \mathbf{x}')) \quad (11.23)$$

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \exp(\mathbf{k}(\mathbf{x}, \mathbf{x}')) \quad (11.24)$$

Where $f : \mathcal{X} \mapsto \mathbb{R}$ a real valued function

$\phi : \mathcal{X} \mapsto \mathbb{R}^e$ the explicit mapping

p a polynomial with pos. coefficients

\mathbf{k}_3 a Kernel over $\mathbb{R}^e \times \mathbb{R}^e$

Proofs

Proof 11.1: Property 11.3 The kernel matrix is positive-semidefinite:

Let $\phi : \mathcal{X} \mapsto \mathbb{R}^d$ and $\Phi = [\phi(\mathbf{x}_1) \dots \phi(\mathbf{x}_n)]^\top \in \mathbb{R}^{d \times n}$.

Thus: $\mathbf{k} = \Phi^\top \Phi \in \mathbb{R}^{n \times n}$.

$$\mathbf{v}^\top \mathbf{k} \mathbf{v} = \mathbf{v}^\top \Phi^\top \Phi \mathbf{v} = (\Phi \mathbf{v})^\top \Phi \mathbf{v} = \|\Phi \mathbf{v}\|_2^2 \geq 0$$

Definition 11.13 Gaussian/Squared Exp. Kernel/ Radial Basis Functions (RBF):

Is an infinite dimensional smooth kernel $\mathbf{k} \in \mathcal{C}^{\infty}$ with some useful properties

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\theta^2}\right) \approx \begin{cases} 1 & \text{if } \mathbf{x} \text{ and } \mathbf{y} \text{ close} \\ 0 & \text{if } \mathbf{x} \text{ and } \mathbf{y} \text{ far away} \end{cases} \quad (11.16)$$

Explanation 11.1 (Threshold θ). $2\theta \in \mathbb{R}$ corresponds to a threshold that determines how close input values need to be in order to be considered similar:

$$\mathbf{k} = \exp\left(-\frac{\text{dist}^2}{2\theta^2}\right) \approx \begin{cases} 1 & \text{if dist} \ll \theta \\ 0 & \text{if dist} \gg \theta \end{cases}$$

or in other words how much we believe in our data i.e. for smaller length scale we do trust our data less and the admirable functions vary much more.

Examples

Example 11.1 Calculating the Kernel by hand:

Let : $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ $\phi(\mathbf{x}) \mapsto \{x_1^2, x_2^2, \sqrt{2}x_1, x_2\}$
 $\phi : \mathbb{R}^{d=2} \mapsto \mathbb{R}^{D=3}$

We can now have a decision boundary in this 3-D feature space \mathcal{V} of ϕ as:

$$\begin{aligned} \beta_0 + \beta_1 x_1^2 + \beta_2 x_2^2 + \beta_3 \sqrt{2}x_1 x_2 &= 0 \\ \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle &= \left\langle \{x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1}, x_{i2}\}, \{x_{j1}^2, x_{j2}^2, \sqrt{2}x_{j1}, x_{j2}\} \right\rangle \\ &= x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2x_{i1} x_{i2} x_{j1} x_{j2} \end{aligned}$$

Operation Count:

- $2 \cdot 3$ operations to map \mathbf{x}_i and \mathbf{x}_j into the 3D space \mathcal{V} .
- Calculating an inner product of $\langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$ with 3 additional operations.

Example 11.2 Calculating the Kernel using the Kernel Trick:

$$\begin{aligned} \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle^2 = \langle \{x_{i1}, x_{i2}\}, \{x_{j1}, x_{j2}\} \rangle^2 \\ &:= \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j) \\ &= (x_{i1} x_{j1} + x_{i2} x_{j2})^2 \\ &= x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2x_{i1} x_{i2} x_{j1} x_{j2} \end{aligned}$$

Operation Count:

- 2 multiplications of $\mathbf{x}_{i1} \mathbf{x}_{j1}$ and $\mathbf{x}_{i2} \mathbf{x}_{j2}$.
- 1 operation for taking the square of a scalar.

Conclusion The Kernel trick needed only 3 in comparison to 9 operations.

Example 11.3 Stationary Kernels:

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \exp\left(\frac{(\mathbf{x} - \mathbf{y})^\top \mathbf{M} (\mathbf{x} - \mathbf{y})}{h^2}\right)$$

is a stationary but not an isotropic kernel.

Time Series

State Space Models

Definition 12.1 State Variables \mathbf{x} :

Is the smallest set of variables $\{x_1, \dots, x_n\}$ that are fully capable of describing the state of our system which is usually *hidden* and not directly observable.

Definition 12.2 State Space \mathcal{X} :

Is the n -dimensional space spanned by the state variables??:
 $\mathbf{x} = [x_1 \dots x_n]^T \in \mathcal{S} \subseteq \mathbb{R}^n$ (12.1)

Definition 12.3

Input/Control Variables $\mathbf{u} \in \mathcal{A}$:

Are a variables \mathbf{u} of the *transition model*^[def. 12.5] that influence the propagation of to the state variables \mathbf{x} .

Definition 12.4

Output/Measurement Variables/State Observations:

Are a variables \mathbf{y} that are directly related to the state space \mathbf{x} and are usually observable by us.

Definition 12.5 Transition Model f :

Describes the transition of the state \mathbf{x} over time.

Definition 12.6

Measurment/Output/Observation Model h :

Describes the mapping of the state \mathbf{x} onto the output \mathbf{y} .

Definition 12.7 (Discrete) State Space Model:

$$\mathbf{x}^{k+1} = f(t, \mathbf{x}^k, \mathbf{u}^k) \quad t = 1, \dots, K \quad (12.2)$$

$$\mathbf{y}^k = h(t, \mathbf{x}^k, \mathbf{u}^k) \quad (12.3)$$

Markov Models

Definition 13.1 States $\mathcal{S} = \{s_1, \dots, s_n\}$:

A state s_i encodes all information of the current configuration of a system.

Definition 13.2

Markovian Property/Memorylessness:

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space with a filtration $(\mathcal{F}_s, s \in I)$, for some index set^[def. 51.1]; and let (S, \mathcal{S}) be a measurable space^[def. 65.7].

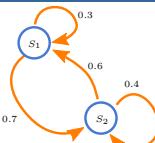
A (S, \mathcal{S}) -valued stochastic process $X = \{X_t : \Omega \rightarrow S\}_{t \in I}$ adapted to the filtration is said to possess the Markov property if:

$$\mathbb{P}(X_t \in A | \mathcal{F}_s) = \mathbb{P}(X_t \in A | X_s) \quad \forall A \in \mathcal{S} \quad s, t \in I \quad \text{s.t. } s < t \quad (13.1)$$

1. Markov Chains

Definition 13.3 Markov Chain:

Is a sequence of random variables $\{X_i\}_{i \in \mathcal{T}}$ ^[def. 69.3] that processes the markovian property^[def. 13.2] i.e. each state X_t depend only on the previous state X_{t-1} :



$\mathbb{P}(X_t = x | X_{t-1} = x_{t-1}, \dots, X_1 = x_1) = \mathbb{P}(X_t = x | X_{t-1} = x_{t-1})$

Definition 13.4 Initial Distribution q_0 :

Describes the initial distribution of states:

$$q_0(s_i) = \mathbb{P}(X_0 = s_i) \quad \forall s_i \in \mathcal{S}$$

$$\iff q_0 = [q_0(s_1) \dots q_0(s_n)] \quad (13.2)$$

Definition 13.5 Transition Probability $p_{ij}(t)$:
 Is the probability of a random variable X_t in state s_i to transition into state s_j :

$$p_{ij}(t) = \mathbb{P}(X_{t+1} = s_j | X_t = s_i) \quad \forall s_i, s_j \in \mathcal{S} \quad (13.3)$$

Definition 13.6 n th Transition Probability $p_{ij}^{(n)}(t)$:
 denotes the probability of reaching state s_j from state s_i in n steps:

$$p_{ij}^{(n)}(t) = \mathbb{P}(X_{t+n} = s_j | X_t = s_i) \quad \forall s_i, s_j \in \mathcal{S} \quad (13.4)$$

Definition 13.7 Transition Matrix $P(t)$:

The transition probabilities eq. (13.4) can be represented by a *row-stochastic matrix*^{??} $P(t)$ where the i^{th} row represents the transition probabilities for the i^{th} state s_i i.e.

From	To j	
0.3	0.7	
0.4		0.6

Corollary 13.1 Row stochastic matrices and Graphs:

Row stochastic matrices?? represent graphs where the outgoing edges must sum to one:

$$\sum \delta^+(s_i) = 1 \quad (13.5)$$

1. Simulating Markov Chains

Corollary 13.2

proof 13.1

Realization of a Markov Chain:

$$\mathbb{P}(X_0 = x_0, \dots, X_N = x_N) = q_0(x_1) \sum_{n=1}^N p_{n-1,n}(t)$$

Algorithm 13.1 Forward Sampling:

Input: $q(\mathbf{x}_0)$ and P
 Output: $\mathbb{P}(X_{0:N})$
 Sample $x_0 \sim \mathbb{P}(X_0)$
 for $j = 1, \dots, n$ do

$$x_j \sim \mathbb{P}(X_j | X_{j-1} = x_{j-1})$$

 5: end for

2. State Distributions

Definition 13.8

Probability Distribution of the States q_{n+1} :

$$\begin{aligned} q_{n+1}(s_j) &= \mathbb{P}(X_{n+1} = s_j) \quad \forall s_i \in \mathcal{S} \\ &= \sum_{i=1}^n \mathbb{P}(X_n = s_i) \mathbb{P}(X_{n+1} = s_j | X_n = s_i) \\ &= \sum_{i=1}^n q_n(s_i) p_{i,j}(t) \end{aligned} \quad (13.6)$$

$$\begin{aligned} q_{n+1} &= [q_{n+1}(s_1) \dots q_{n+1}(s_n)] \\ &= q_n P(t) \\ &= [q_n(s_1) \dots q_n(s_n)] \begin{bmatrix} p_{1,1} & p_{1,2} & \dots & p_{1,n} \\ p_{2,1} & p_{2,2} & \dots & p_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ p_{n,1} & p_{n,2} & \dots & p_{n,n} \end{bmatrix} (t) \end{aligned}$$

Corollary 13.3 Time-homogeneous Markov Transition Probabilities:

$$q_{n+1} = q_0 P^{n+1} \quad [\text{proof 13.2}]$$

Definition 13.9 Stationary Distribution:

A markov chain has a stationary distribution if it satisfies:

$$\lim_{N \rightarrow \infty} q_N(s_i) = \lim_{N \rightarrow \infty} \mathbb{P}(X_N = s_i) = \pi_i \quad \forall s_i \in \mathcal{S}$$

$$\lim_{N \rightarrow \infty} q_N = [\pi_1 \dots \pi_n] \iff q = \mathbb{P}(N) \quad (13.8)$$

Corollary 13.4 Existence of Stationary Distributions:
 A Markov Chain has a stationary distribution if and only if at least one state is *positive recurrent*!

3. Properties of States

Definition 13.10 Absorbing State/Sink:

Is a state s_i that once entered cannot be left anymore:

$$p_{ij}^{(n)}(t) = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases} \quad (13.9)$$

Definition 13.11 Accessible State $s_i \rightarrow s_j$:

A state s_j is accessible from state s_i iff:

$$\exists n : \mathbb{P}_{ij}^{(n)}(t) > 0 \quad (13.10)$$

Definition 13.12 Communicating States $s_i \leftrightarrow s_j$:

Two states s_j and s_i are communicating iff:

$$\exists n_1 : \mathbb{P}_{ij}^{(n_1)}(t) > 0 \quad \wedge \quad \exists n_2 : \mathbb{P}_{ji}^{(n_2)}(t) > 0 \quad (13.11)$$

Definition 13.13 Periodicity of States: A state s_i has period k if any return to state s_i must occur in multiples of k time steps.

In other words k is the *greatest common divisor* of the number of transitions by which state s_i can be reached, starting from itself:

$$k = \gcd\{n > 0 : \mathbb{P}_{ii}^{(n)} = \mathbb{P}(X_n = s_i | X_0 = s_i) > 0\} \quad (13.12)$$

Definition 13.14 Aperiodic State $k = 1$:
 Is a state s_i with periodicity^[def. 13.13] of one $\Leftrightarrow k = 1$

Corollary 13.5 : A state s_i is aperiodic if there exist two consecutive numbers k and $k+1$ s.t. the chain can be in state s_i at both time steps k and $k+1$.

Corollary 13.6 Absorbing State: An absorbing state is an aperiodic state.

Explanation 13.1 (Definition 13.14). Returns to state s_i can occur at irregular times i.e. the state is not predictable.
 In other words we cannot predict if the state will be revisited in multiples of k times.

4. Characteristics of Markov Processes/Chains

Definition 13.15

Time-homogeneous/Stationary Markov Chain:

are markov chains^[def. 13.3] where the transition probability is independent of time:

$$\mathbb{P}_{ij} = \mathbb{P}(X_t = s_j | X_{t-1} = s_i) = \mathbb{P}(X_{t-\tau} = s_j | X_{t-\tau} = s_i) \quad \forall \tau \in \mathbb{N}_0 \quad (13.13)$$

Corollary 13.7 Transition Matrices of Stationary MCs:

Transition matrices of time-homogeneous markov chain are constant/time independent:

$$P(t) = P \quad (13.14)$$

Definition 13.16 Aperiodic Markov Chain: Is a markov chain where all states are aperiodic:

$$\gcd\{n > 0 : \mathbb{P}_{ii}^{(n)} = \mathbb{P}(X_n = s_i | X_0 = s_i) > 0\} = 1 \quad \forall i \in \{1, \dots, n\} \quad (13.15)$$

Definition 13.17 Irreducible Markov Chain: Is a Markov chain that has only *communicating states*^[def. 13.12]:

$$s_j \leftrightarrow s_i \quad \forall i, j \in \{1, \dots, n\} \quad (13.16)$$

• \Rightarrow no sinks^[def. 13.10]

• \Rightarrow every state can be reached from every other state

Corollary 13.8 : An *irreducible*^[def. 13.17] markov chain is automatically *aperiodic*^[def. 13.16] if it has at least one aperiodic state^[def. 13.14] \Leftrightarrow *ergodic*^[def. 13.18].

Corollary 13.9 : A markov chain is *not-irreducible* if there exist two states with different periods.

Definition 13.18

[example 13.1]

Ergodic Markov Chain: A finite markov chain is ergodic if there exist some number N s.t. any state s_j can be reached from any other state s_i in any number of steps less or equal to a N .

\Rightarrow a markov chains is ergodic if it is:

- ① *Irreducible*^[def. 13.17]
- ② *Aperiodic*^[def. 13.16]

Corollary 13.10 Stationary Distribution: An ergodic markov chain has a *unique* stationary distribution^[def. 13.9] and converges to it starting from any initial state $q_0(s_i)$

5. Types of Markov Chains

	Observable	Unobservable
Uncontrolled	MC ^[def. 13.3]	HMM ^[def. 14.1]
Controlled	MDP ^[def. 15.1]	POMDP ^[def. 16.1]

6. Markov Chain Monte Carlo (MCMC)

2. Proofs

Proof 13.1: [cor. 13.2]

$$\begin{aligned} \mathbb{P}(X_0 = x_0, \dots, X_N = x_N) &= \mathbb{P}(X_0 = x_0) \cdot \\ &\cdot \mathbb{P}(X_1 = x_1 | X_0 = x_0) \cdot \mathbb{P}(X_2 = x_2 | X_1 = x_1, X_0 = x_0) \cdot \\ &\cdots \mathbb{P}(X_N = x_N | X_{N-1} = x_{N-1}, \dots, X_0 = x_0) \end{aligned}$$

and then simply use the Markovian property

Proof 13.2: Corollary 13.3

$$q_{n+1} = \mathbb{P} q_n = (\mathbb{P} q_n) P = q_0 P^{n+1}$$

3. Examples

Example 13.1 Ergodic Markov Chain:

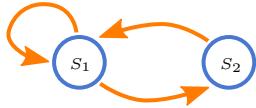


Figure 9: Ergodic for $N = 2$ (can reach s_2 at any $t \leq N$ after $N = 2$)

Hidden Markov Model (HMM)

Definition 14.1 Hidden Markov Model (HMM):

Is a Markov Chain^[def. 13.3] with hidden/latent states S_j that are only partially observable by noisy/indirect observations^[def. 14.2]: It is characterized by the 5-tuple of:

- ① States^[def. 13.1] $S = \{s_1, \dots, s_n\}$
- ② Actions^[def. 15.2] $\mathcal{A}/\mathcal{A}_{s_j} = \{a_1, \dots, a_m\}$
- ③ Observations^[def. 14.2] $\mathcal{O}/\mathcal{O}_{s_j} = \{o_1, \dots, o_m\}$
- ④ Transition Probabilities^[def. 13.5] $P(s_i, s_j)$
- ⑤ Emission/Output Probabilities^[def. 14.3] $e_{ij}(t)$

Definition 14.2 Observations $\mathcal{O} = \{o_1, \dots, o_l\}$: Are indirect or noisy observations that are related to the true states s_j .

Definition 14.3 Emission/ Output Probabilities $e_{ij}(t)$:

Given a state $X_t = s_i$ the output probability is the probability of the output random variable Y_t to be in state o_j :

$$e_{ij}(t) = \mathbb{P}(Y_t = o_j | X_t = s_i) \quad \begin{cases} \forall o_i \in \mathcal{O} \\ \forall s_j \in \mathcal{S} \end{cases} \quad (14.1)$$

Latent $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_T \rightarrow X_{t+1}$

Observed $Y_1 \rightarrow Y_2 \rightarrow \dots \rightarrow Y_T \rightarrow Y_{t+1}$

$e_{ik}(t)$

Markov Decision Processes (MDP)

Definition 15.1 $(\mathcal{S}, \mathcal{A}, \mathbb{P}_a, R_a)$

Markov Decision Process (MDP): A markov decision process is a *controlled* markov process/chain with an associated reward, where the transition can be steered by an actions. It is characterized by the 4-tuple of:

- ① States^[def. 13.1] $\mathcal{S} = \{s_1, \dots, s_n\}$
- ② Actions^[def. 15.2] $\mathcal{A}/\mathcal{A}_{s_j} = \{a_1, \dots, a_m\}$
- ③ Transition Probabilities^[def. 15.3] $\mathbb{P}_a(s_i, s_j)$
- ④ Rewards^[def. 15.4] $r_a(s_i, s_j)$

Definition 15.2

$$\mathcal{A}_{s_i} = \{a_1, \dots, a_m\}$$

Is the set of possible actions from which we can choose at each state and may depend on the state s_j itself.

Definition 15.3 Transition Probability $\mathbb{P}_a(s_j, s_i)(t)$:

is the probability of a random variable X_t in state s_i to transition into state s_j and depends also on the current action a :

$$\mathbb{P}_a(s_j, s_i) = p(s_j|s_i, a) = P(x_{t+1} = s_j | x_t = s_i, a_t = a) \quad \forall s_i, s_j \in \mathcal{S}, \forall a \in \mathcal{A} \quad (15.1)$$

Definition 15.4 Reward

$$r_a(s_i, s_j)$$

is a function or probability distribution that measures the immediate reward and may depend on a any subset of (x_{t+1}, x_t, a) :

$$(x_{t+1}, x_t, a) \mapsto R_{t+1} \in \mathcal{R} \subset \mathbb{R} \quad (15.2)$$

Markov decision processes require us to plan ahead. This is because the immediate reward^[def. 15.4], that we obtain by greedily picking the best action may result in non-optimal local actions.

1. Policies and Values

Definition 15.5

Optimizing Agent / Decision Making Policy $\pi(s_i)$: Is a policy on how to choose an action $a \in \mathcal{A}$ based on a objective/value function^[def. 15.8] and can be deterministic or randomized:

$$\pi : \mathcal{S} \mapsto \mathcal{A} \quad \text{or} \quad \pi : \mathcal{S} \mapsto \mathbb{P}(\mathcal{A}) \quad (15.3)$$

Definition 15.6 Discounting Factor

γ : Is a factor $\gamma \in [0, 1]$ that signifies that future rewards are less valuable then current rewards.

Explanation 15.1 (Definition 15.6). *The reason for the discounting factor is that we may for example not even survive long enough to obtain future payoffs.*

Definition 15.7 Expected Discounted Value $J(\pi)$: Is the discounted expected (reward) of the whole markov process:

$$J(\pi) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \right] \quad (15.4)$$

Definition 15.8 Value Function $V^\pi(x)$: Is the discounted expected reward^[def. 15.4] of the whole markov process given an initial state $X_0 = x$:

$$V^\pi(x) = J(\pi|X_0 = x) \quad (15.5)$$

$$= \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \quad (15.6)$$

$$(15.7)$$

1. Calculating the value of V^π

Definition 15.9

Value Iteration:

$$\begin{aligned} V^\pi(x) &= J(\pi|X_0 = x) \\ &= \mathbb{E}_{x'|x, \pi(x)} [r(x, \pi(x)) + \gamma V^\pi(x')] \\ &= r(x, \pi(x)) + \gamma \mathbb{E}_{x'|x, \pi(x)} [V^\pi(x')] \\ &= r(x, \pi(x)) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x'|x, \pi(x)) V^\pi(x') \end{aligned} \quad [proof 16.1] \quad (15.8)$$

We can now write this for all possible initial states as:

$$V^\pi = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi V^\pi \iff (\mathbf{I} - \gamma \mathbf{P}^\pi) V^\pi = \mathbf{r}^\pi \quad (15.9)$$

with:

$$\begin{aligned} \mathbf{V}^\pi &= \begin{bmatrix} V^\pi(s_1) \\ \vdots \\ V^\pi(s_n) \end{bmatrix} & \mathbf{r}^\pi &= \begin{bmatrix} r^\pi(s_1, \pi(s_1)) \\ \vdots \\ r^\pi(s_n, \pi(s_n)) \end{bmatrix} \\ \mathbf{P}^\pi &= \begin{bmatrix} \mathbb{P}(s_1|s_1, \pi(s_1)) & \mathbb{P}(s_2|s_1, \pi(s_1)) & \dots & \mathbb{P}(s_n|s_1, \pi(s_1)) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{P}(s_1|s_2, \pi(s_2)) & \mathbb{P}(s_2|s_2, \pi(s_2)) & \dots & \mathbb{P}(s_n|s_2, \pi(s_2)) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{P}(s_1|s_n, \pi(s_n)) & \mathbb{P}(s_2|s_n, \pi(s_n)) & \dots & \mathbb{P}(s_n|s_n, \pi(s_n)) \end{bmatrix} \end{aligned}$$

1.1.1. Direct Methods

Corollary 15.1 LU-decomposition

$\mathcal{O}(n^3)$:

The linear system from eq. (15.9): $(\mathbf{I} - \gamma \mathbf{P}^\pi) \mathbf{V}^\pi = \mathbf{r}^\pi$ (15.10)

can be solved directly using Gaussian elimination in polynomial time $\mathcal{O}(n^3)$.

Note – invertibility

If $\gamma < 1$ then $(\mathbf{I} - \gamma \mathbf{P}^\pi)$ is full-rank/invertible as EVs(\mathbf{P}^π) ≤ 1 .

1.1.2. Fixed Point Iteration

Corollary 15.2 Fixed-Point Iteration

$\mathcal{O}(n \cdot |\mathcal{S}|)$:

The linear system from eq. (15.9) can be solve using fixed-point iteration^[def. 6.2, 30] in at most $\mathcal{O}(n \cdot |\mathcal{S}|)$ (if every state s_i is connected to every other state $s_j \in \mathcal{S}$)

Algorithm 15.1 Fixed Point Iteration:

Input: Initial Guess: V_0^π i.e. 0

1: **for** $t = 1, \dots, T$ **do**
2: Use the fixed point method:
$$V_t^\pi = \phi V_t^\pi = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi V_{t-1}^\pi \quad (15.11)$$

3: **end for**

Corollary 15.3

Policy Iteration Contraction [proof 16.2]:

Fixed point iteration of policy iteration is a contraction^[def. 59, 63] that leads to a fixed point V^π with a rate depending on the discount factor γ .

$$\|V_t^\pi - V^\pi\| = \|\phi V_t^\pi - \phi V^\pi\| \leq \gamma \|V_{t-1}^\pi - V^\pi\| = \gamma^t \|V_0^\pi - V^\pi\| \quad (15.12)$$

Explanation 15.2.

- $\gamma \downarrow$: the less we plan ahead/the smaller we choose γ the shorter it takes to converge. But on the other hand we only care greedily about local optima and might miss global optima.
- $\gamma \uparrow$: the more we plan ahead/the larger we choose γ the longer it takes to converge but we will explore all possibilities. But for to large γ we will simply keep exploring without sticking to a optimal point

Note contraction

For a contraction:

- A unique fixed point exists
- We converge to the fixpoint

2. Choosing The Policy

Question how should we choose the π ? **Idea** compute $J(\pi)$ for every possible policy:

$$\pi^* = \arg \max \mathbb{J}(\pi) \quad (15.13)$$

Problem this is unfortunately infeasible as there exist $m^n = |\mathcal{A}|^{|S|}$ policies that we need to calculate the value for.

Note

The problem is that J/V^π depend on π but if we do not know π yet we cannot compute those.

1.2.1. Greedy Policy

Definition 15.10 Greedy Policy:

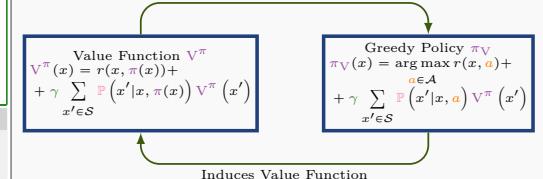
Assuming we know V^{*t-1} then we could choose a greedy policy:

$$\begin{aligned} a^* &= \pi_t(x) \\ &:= \arg \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x'|x, a) V^{*t-1}(x') \end{aligned} \quad (15.14)$$

- Given a policy π however we can calculate a value function V^π

- Given a value function V we can induce a greedy policy^[def. 15.10] π w.r.t. V

Induces Policy



Induces Value Function

Theorem 15.1 Optimality of Policies [Bellman]:

A policy π_V is optimal if and only if it is greedy w.r.t. its induced value function

Definition 15.11 Non-linear Bellman Equation: States that the optimal value is given by the action/policy that maximizes the value function eq. (15.8):

$$V^*(x) = \max_{a \in \mathcal{A}} \left[r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x'|x, a) V^*(x') \right] \quad (15.15)$$

$$:= \max_{a \in \mathcal{A}} Q^*(x, a) \quad (15.16)$$

Note

This equation is non-linear due to the max in comparison to eq. (15.8).

1.2.2. Policy Iteration

Algorithm 15.2 Policy Iteration:

Initialize: Random Policy: π

1: **while** Not converged $t = t + 1$ **do**

2: Compute $V^{*t}(x)$
$$V^{*t}(x) = r(x, \pi(x)) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x'|x, \pi(x)) V^{*t-1}(x')$$

3: Compute greedy policy π_G :

$$\pi_G(x) = \arg \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x'|x, a) V^{*t-1}(x')$$

4: Set $\pi_{t+1} \leftarrow \pi_G$

5: **end while**

Algorithm 15.2

Pros

- Monotonically improves $V^{*t} \geq V^{*t-1}$
- is guaranteed to converge to an optimal policy/solution π^* in polynomial #iterations: $\mathcal{O}\left(\frac{n^2 m}{1-\gamma}\right)$

Cons

- Complexity per iteration requires to evaluate the policy V^π which requires us to solve a linear system.

1.2.3. Value Iteration

Definition 15.12 Value to Go

$V_t(x)$: Is the maximal expected reward if we start in state x and have t time steps to go.

Algorithm 15.3 Value Iteration [proof 16.3]:

Initialize: $V_0(x) = \max_{a \in \mathcal{A}} r(x, a)$

```

1: for  $t = 1, \dots, \infty$  do
2:   Compute:
3:    $Q_t(x, a) = r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x'|x, a) V_{t-1}(x')$   $\forall a \in \mathcal{A}$ 
4:   for all  $x \in \mathcal{S}$  let:
       $V_t(x) = \max_{a \in \mathcal{A}} Q_t(x, a)$ 
5:   if  $\max_{x \in \mathcal{S}} |V_t(x) - V_{t-1}(x)| \leq \epsilon$  then
6:     break
7:   end if
8: end for
9: Choose greedy policy  $\pi_{V_t}$  w.r.t.  $V_t$ 
```

Corollary 15.4

Value Iteration Contraction:

Algorithm 15.3 is guaranteed to converge to a ϵ optimal policy:

$$\begin{aligned} \|V_t - V^*\|_\infty &\leq \gamma^t \|V_0 - V^*\|_\infty \\ \Rightarrow t &\approx \ln \frac{\epsilon}{\gamma} \|V_0 - V^*\|_\infty \quad \text{for } \|V_t - V^*\|_\infty \leq \epsilon \end{aligned} \quad (15.17)$$

Algorithm 15.3

Pros

- Finds ϵ -optimal solution in polynomial #iterations $\mathcal{O}(\ln \frac{1}{\epsilon})$ [cor. 15.4].
- Complexity per iteration requires us to solve a linear system $\mathcal{O}(m \cdot n \cdot s) = \mathcal{O}(|\mathcal{A}| \cdot |\mathcal{S}| \cdot s)$ where s is the number of states we can reach. For small s and small m we are roughly linear w.r.t. the states $\mathcal{O}(n) = \mathcal{O}(|\mathcal{S}|)$

Cons

- Only ϵ -optimal solution.

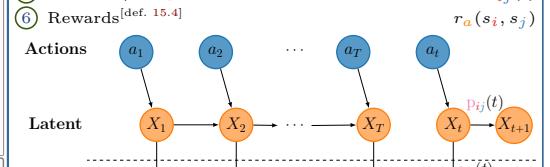
Partially Observable MDP (POMDP)

Definition 16.1 $(\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathbb{P}_a, E, R_a)$

Partially Observable Markov Decision Process:

A (POMDP) is a markov decision process^[def. 15.1] with hidden markov states^[def. 14.1]. It is characterized by the 6-tuple of:

- States^[def. 13.1] $\mathcal{S} = \{s_1, \dots, s_n\}$
- Actions^[def. 15.2] $\mathcal{A}/\mathcal{A}_{s_j} = \{a_1, \dots, a_m\}$
- Observations^[def. 14.2] $\mathcal{O}/\mathcal{O}_{s_j} = \{o_1, \dots, o_m\}$
- Transition Probabilities^[def. 15.3] $\mathbb{P}_a(s_i, s_j)$
- Emission/Output Probabilities^[def. 14.3] $e_{ij}(t)$
- Rewards^[def. 15.4] $r_a(s_i, s_j)$



Explanation 16.1.

Now our agent has only some indirect noisy observation of true state.

1. POMDPs as MDPs

POMDPs can be converted into belief state?? MDPs^[def. 15.1] by introducing a belief state space \mathcal{B} .

Definition 16.2 History H_t : Is a sequence of actions, observations and rewards:
 $H_t = \{\langle a_0, o_0, r_0 \rangle, \dots, \langle a_0, o_0, r_0 \rangle\}$

Definition 16.3 Belief State Space \mathcal{B} : Is a $|\mathcal{S}| - 1$ dimensional simplex or ($|\mathcal{S}|$ -dimensional probability vector^[def. 59.71]) whose elements b are probabilities:

$$\mathcal{B} = \Delta(|\mathcal{S}|) = \left\{ b_t \in [0, 1]^{|\mathcal{S}|} \mid \sum_{x=1}^n b_t(x) = 1 \right\} \quad (16.1)$$

Definition 16.4 Belief State $b_t \in \mathcal{B}$: Is a probability distribution over the states \mathcal{S} conditioned on the history H_t ^[def. 16.2].

1. Transition Model

Definition 16.5 POMDP State/Posterior Update: [proof 16.5]

$$\begin{aligned} b_{t+1}(s_i) &= \mathbb{P}(X_{t+1} = s_i | Y_{t+1} = o_k) \\ &= \frac{1}{Z} \mathbb{P}(Y_{t+1} = o_k | X_{t+1} = s_i, a_t) \\ &\quad \cdot \sum_{x' \in \mathcal{S}} b_t(s_j) \mathbb{P}(X_{t+1} = s_i | X_t = s_j, a_t) \end{aligned} \quad (16.2)$$

Definition 16.6 Stochastic Observation Model:

$$\begin{aligned} \mathbb{P}(Y_{t+1} = o_k | b_t, a_t) &= \sum_{s_i \in \mathcal{S}} b_t(s_i) \mathbb{P}(Y_{t+1} = o_k | X_t = s_i, a_t) \\ &\quad \text{eq. (13.13)} \end{aligned} \quad (16.3)$$

2. Reward Function

Definition 16.7 POMDP Reward Function: $r(b_t, a_t) = \sum_{s_j \in \mathcal{S}} b_t(s_j) r(s_j, a_t)$ (16.4)

Note

For finite horizon T , the set of reachable belief states is finite however exponential in T .

2. Proofs

1. Markov Decision Processes

Proof 16.1: [def. 15.8]

$$\begin{aligned} \mathbb{V}^\pi(x) &= \mathbb{E}_{X_{1:\infty}} \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \\ &= \mathbb{E}_X \left[\gamma^0 r(X_0, \pi(X_0)) + \sum_{t=1}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \\ &\stackrel{\gamma^0=1}{=} r(x, \pi(x)) + \mathbb{E}_X \left[\sum_{t=1}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \\ &\stackrel{\text{re-index}}{=} r(x, \pi(x)) + \mathbb{E}_X \left[\sum_{t=0}^{\infty} \gamma^{t+1} r(X_{t+1}, \pi(X_{t+1})) \mid X_0 = x \right] \\ &= r(x, \pi(x)) + \gamma \mathbb{E}_X \left[\sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_0 = x \right] \\ &= r(x, \pi(x)) + \gamma \mathbb{E}_{X_1} \left[\mathbb{E}_{X_{2:\infty}} \left[\sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_1 = x' \right] \mid X_0 = x \right] \\ &\stackrel{\text{law 65.7}}{=} r(x, \pi(x)) \\ &\quad + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, \pi(x)) \mathbb{E}_{X_{2:\infty}} \left[\sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_1 = x' \right] \\ &\stackrel{\text{eq. (13.13)}}{=} r(x, \pi(x)) \\ &\quad + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, \pi(x)) \mathbb{E}_{X_{2:\infty}} \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x' \right] \end{aligned}$$

Proof 16.2 [cor. 15.3]: Consider $\mathbb{V}, \mathbb{V}' \in \mathbb{R}^n$ and let ϕ :

$$\phi x := r^\pi + \gamma \mathbb{P}^\pi x \implies \phi \mathbb{V}^\pi = \mathbb{V}^\pi$$

then it follows:

$$\begin{aligned} \|\phi \mathbb{V} - \phi \mathbb{V}'\| &= \left\| \cancel{r^\pi} + \gamma \mathbb{P}^\pi \mathbb{V} - \cancel{r^\pi} - \gamma \mathbb{P}^\pi \mathbb{V}' \right\| \\ &= \left\| \gamma \mathbb{P}^\pi (\mathbb{V} - \mathbb{V}') \right\| \\ &\stackrel{\text{eq. (59.91)}}{\leq} \gamma \left\| \mathbb{P}^\pi \right\| \cdot \left\| (\mathbb{V} - \mathbb{V}') \right\| \\ &\stackrel{\text{i.e. L2}}{\leq} \gamma \cdot 1 \cdot \left\| (\mathbb{V} - \mathbb{V}') \right\|_2 \end{aligned}$$

Proof 16.3: algorithm 15.3

$$\begin{aligned} \mathbb{V}_0(x) &= \max_{a \in \mathcal{A}} r(x, a) \\ \mathbb{V}_1(x) &= \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, a) \mathbb{V}_0(x') \\ \mathbb{V}_{t+1}(x) &= \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, a) \mathbb{V}_t(x') \end{aligned}$$

Proof 16.4: [cor. 15.4] Let $\phi : \mathbb{R}^n \mapsto \mathbb{R}^n$, with:

$$\begin{aligned} (\phi \mathbb{V}^*) (x) &= Q(x, a) = \max_a \left[r(x, a) + \gamma \sum_{x'} \mathbb{P}(x' | x, a) \right] \\ \text{Bellman's theorem 15.1} \quad \phi \mathbb{V}^* &= \mathbb{V}^* \\ \|\phi \mathbb{V} - \phi \mathbb{V}'\|_\infty &= \max_x |(\phi \mathbb{V})(x) - (\phi \mathbb{V}')(x)| \\ &= \max_x \left| \max_a Q(x, a) - \max_{a'} Q'(x, a') \right| \\ &\stackrel{\text{Property 54.9}}{\leq} \max_x \max_a |Q(x, a) - Q'(x, a)| \\ &= \max_{x, a} \left| \frac{1}{\gamma} + \gamma \sum_{x'} \mathbb{P}(x' | x, a) \mathbb{V}(x') - \frac{1}{\gamma} - \gamma \sum_{x'} \mathbb{P}(x' | x, a) \mathbb{V}'(x') \right| \\ &= \gamma \max_{x, a} \left| \sum_{x'} \mathbb{P}(x' | x, a) (\mathbb{V}(x') - \mathbb{V}'(x')) \right| \\ &\stackrel{\leq 1}{\leq} \\ &\stackrel{\text{eq. (59.91)}}{\leq} \gamma \max_{x, a} \left| \sum_{x'} \mathbb{P}(x' | x, a) \right| \cdot |(\mathbb{V}(x') - \mathbb{V}'(x'))| \\ &\leq \gamma \cdot 1 \cdot \|(\mathbb{V}(x') - \mathbb{V}'(x'))\|_\infty \end{aligned}$$

Note

For the policy iteration the calculation was easier as the rewards canceled, however here we have the max.

2. MDPs

Proof 16.5: Defintion 16.5 Directly by definition 7.5 and its corresponding proof 10.4 with additional action a_t :

$$\begin{aligned} b_{t+1}(s_i) &= \mathbb{P}(X_{t+1} = s_i | y_{1:t+1}) \\ &= \frac{1}{Z} \mathbb{P}(y_{1:t+1} | s_i) \sum_{j=1}^Z \underbrace{\mathbb{P}(X_t = s_j | y_{1:t})}_{b_t(s_j)} \mathbb{P}(s_i | s_j) \end{aligned}$$

Reinforcement Learning

Now we are working with an *unknown* MDP^[def. 15.1] meaning that:

- ① we do no longer know the transition model^[def. 15.3]
- ② We do no longer know the reward function
- ③ We might not even know all the states
However we can observe them when taking steps.

Note

- Reinforcement learning is different than supervised learning as the data is no longer i.i.d. (data depends on previous action).
- Need to do exploration vs exploitation in order to learn policy and reward functions.

Definition 17.1 Agent:

Is the *learner/decision maker* of our *unknown* MDP.

Definition 17.2 Environment:

Is the representation of the world in which our agents acts.

Definition 17.3 On-Policy Learning:

At any given time the agent has full control which actions to pick.

Definition 17.4 Off-Policy Learning:

The agent has to fix a policy in advance based on behavioral observations.

Definition 17.5 Trajectory

τ : Is a set of consecutive 3-tuples of states, actions and rewards:
 $\tau = \{s_t, a_t, r_t\} \quad t = 1, \dots, \tau$ (17.1)

Definition 17.6 Episodic Learning:

Is a setting where we generate multiple K -episodes of different trajectories $\{\tau^{(k)}\}_{k=1}^K$ from which the agent can learn.

Explanation 17.1.

For each episode the agent starts in a random state and follows a policy.

1. Model Based Reinforcement Learning

Proposition 17.1 Model Based RL:

Try to learn the MDP^[def. 15.1] by:

- ① Estimating
 - the transition probabilities^[def. 15.3] $p_a(s_i, s_j)$
 - the reward function^[def. 15.4] $r(b_t, a_t)$
- ② Optimizing the policy of the estimated MDP

1. Estimating Transitions and Rewards

Formula 17.1 Estimating Transitions and Rewards:

Given a data set $D = \{(x_0, a_0, r_0, x_1), (x_1, a_1, r_1, x_2), \dots\}$ we estimate the transitions and rewards using a categorical distribution^[def. 66.23]:

$$N_{s_i|s_j, a} := \sum_{k=1}^t \delta_{(x_{k+1}=s_i | X_k=s_j, A_k=a)} \quad (17.2)$$

$$N_{s_j, a} := \sum_{k=1}^t \delta_{(X_k=s_j, A_k=a)} \quad (17.3)$$

$$p_a(s_i, s_j) \approx \frac{N_{s_i|s_j, a}}{N_{s_j, a}} \quad (17.4)$$

$$r(s_i, a) \approx \frac{1}{N_{s_i, a}} \sum_{k=1}^t \delta_{(X_k=s_i, A_k=a)} r(X_k, A_k) \quad (17.5)$$

2. Choosing the next step

How should we choose the action $a \in \mathcal{A}$ in order to balance exploration vs exploitation?

3. ϵ_t Greedy Learning

Algorithm 17.1 Epsilon Greedy Learning:

```

1: for  $t = 1, \dots, T$  do
2:   Pick next action
       $a_t = \begin{cases} \arg \max_a Q_t(a) & \text{with probability } \epsilon_t \\ \text{random } a & \text{with probability } 1 - \epsilon_t \end{cases}$ 
3: end for

```

Corollary 17.1 Necessary Condition for Convergence:

If the sequence ϵ_t satisfies the Robbins Monro (RM) conditions

$$\sum_t \epsilon_t < \infty, \quad \sum_t \epsilon_t^2 < \infty \quad (\text{i.e. } \epsilon_t = 1/t) \quad (17.6)$$

then algorithm 17.1 converges to an optimal policy with probability one.

Pros

- Simple
- Clearly sub optimal actions are not eliminated fast enough

4. The R_{\max} Algorithm

Algorithm 17.2 [Brafman & Tennenholtz '02]

R-max Algorithm:

Initialize every state with:

$$\hat{r}(s_t, a) = R_{\max} \quad \hat{p}_a(X_{t+1}|X_t = s_i, a) = 1 \quad (17.7)$$

Set min. number Δ of observations for policy update

Compute Policy π_1 of the MDP^[def. 15.1] using (\hat{p}, \hat{r}) :

$$\pi_t$$

```

1: for  $k = 1, \dots, K$  do
2:   Choose  $a = \pi_t(x_t)$  and observe  $(s, r)$ 
3:   Calculate:
       $N_{x_t, a} += 1 \quad r(x_t, a) += r(x_t, a) \quad (17.8)$ 
       $N_{x_{t+1}|x_t, a} += 1$   $\quad (17.9)$ 
4:   if  $k == \Delta$  then
5:     Re-calculate (based on eqs. (17.4) and (17.5)):
       $\hat{r}(s_t, a) = R_{\max} \quad \hat{p}_a(X_{t+1}|X_t = s_i, a) = 1$ 
      and update the policy  $\pi_t = \pi_t(\hat{p}, \hat{r})$ 
6:   end if
7: end for

```

Note

Other ways of updating the policy at certain times exist.

Problems

- ① **Cons**
 - Memory: for all $a \in \mathcal{A}$, $x_{t+1}, x_t \in \mathcal{X}$ we need to store $\hat{p}_a(x_{t+1}|x_t, a)$ and $\hat{r}(s_t, a)$ which results in $|\mathcal{S}|^2 |\mathcal{A}|$ (for dense MDP).
 - Computation Time: We need to calculate the π_t using policy (?? 1.2.2) or value iteration (?? 1.2.3) $|\mathcal{A}| \cdot |\mathcal{S}|$ whenever we update our policy.

1.4.1. How many transitions do we need?

Proposition 17.2 [proof 17.1]

Number of Samples to bound Reward:

$$\mathbb{P}(\hat{r}(s, a) - r(s, a) \leq \epsilon) \geq 1 - \delta \iff n \in \mathcal{O}\left(\frac{R_{\max}^2 \log \frac{1}{\delta}}{\epsilon^2}\right) \quad (17.10)$$

Theorem 17.1 :

Every T timesteps, with high probability, R_{\max} either:

- Obtain near optimal reward, or
- Visits at least one unknown state-action pair

Theorem 17.2 Performance of R-max:

With probability $\delta - 1$, R_{\max} will reach an ϵ -optimal policy in a number of steps that is polynomial in $|\mathcal{X}|, |\mathcal{A}|, T, 1/\epsilon$.

2. Model Free Reinforcement Learning

Proposition 17.3 Model Free RL:

Tries to estimate the value function^[def. 15.8] directly in order to act greedily upon it.

- Policy Gradient Methods
- Actor Critic Methods

1. Temporal Difference (TD)

Assume we fix a random initial policy π and s.t. we have $\hat{V}_0(s_i)$.

Goal: want to calculate an unknown value function V^π . If the reward and the next states are stochastic variables (R, X) we can calculate the reward using eq. (15.8):

$$\hat{V}^\pi(x_t) = \mathbb{E}_{X_{t+1}, R} [R + \gamma \hat{V}^\pi(X')|X, a] \quad (17.11)$$

Now assume we observe a single example

$$(X_{t+1} = s_j, a, r, X_t = s_i)$$

then we can use monte carlo sampling^[def. 67.6] with a single sample to approximate the expectation ineq. (17.11):

$$\hat{V}_{t+1}^\pi(s_i) = r + \gamma \hat{V}_t^\pi(s_j)$$

Problem: high variance of estimates \Rightarrow average with previous estimate.

Definition 17.7 Temporal Difference (TD) Learning:

$$\hat{V}(x_{t+1}) = (1 - \alpha_t) \hat{V}(x_t) + \alpha_t (r + \gamma \hat{V}(x_{t+1})) \quad (17.12)$$

Corollary 17.2 Necessary Condition for Convergence:

If the learning rate α_t satisfies the Robbins Monro (RM) conditions

$$\sum_t \alpha_t < \infty, \quad \sum_t \alpha_t^2 < \infty \quad (\text{i.e. } \alpha_t = 1/t) \quad (17.13)$$

and all state-action pairs (s_i, a_j) are chosen infinitely often, then we converge to the correct value function:

$$\mathbb{P}(\hat{V} \rightarrow V^\pi) = 1 \quad (17.14)$$

2. Q-Learning

Definition 17.8 Action Value/Q-Function:

$$Q \quad (17.15)$$

2.2.1. Policy Gradients

2.2.2. Actor-Critic Methods

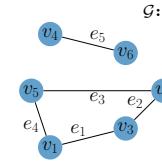
3. Proofs

Proof 17.1: proposition 17.2 using hoeffdings bound^[def. 65.38] with δ and $b - a = R_{\max}$.

Graph Theory

Definition 18.1 Graph

A graph G is a pair $G = (\mathcal{V}, \mathcal{E})$ of a finite set of vertices \mathcal{V} and a multi set of edges \mathcal{E} .



Definition 18.2 Order

$n = |\mathcal{V}|$: The order of a graph is the cardinality of its vertex set.

Definition 18.3 Size

$m = |\mathcal{E}|$: The size of a graph is the number of its edges.

Corollary 18.1 n-Graph:

Is a graph G of order n .

Corollary 18.2 (p, q)-Graph:

Is a graph G of order p and size q .

Vertices

Definition 18.4 Vertices/Nodes

\mathcal{V} : Is a set of entities of a graph connected and related by edges in some way:

Definition 18.5 Neighbourhood

$N(v_i)$: The neighborhood of a vertex $v_i \in \mathcal{V}$ is the set of all adjacent vertices:

$$N(v_i) = \{v_k \in \mathcal{V} : \exists e_k = \{v_i, v_k\} \in \mathcal{E}, \forall v_j \in \mathcal{E}\} \quad (18.1)$$

Degree Matrix

Definition 18.6 Degree of a Vertex

δ : The degree of a vertex v is the cardinality of the neighborhood – the number of adjacent vertices:

$$\deg(v_i) = \delta(v) = |N(v)| = \sum_{j=1}^{j < i} \mathbf{A}_{ij} \quad (18.2)$$

Definition 18.7 Degree Matrix

\mathbf{D} : Given a graph $G = (\mathcal{V}, \mathcal{E})$ its degree matrix is a diagonal matrix $\mathbf{D} \in \mathbb{N}^{n,n}$ defined as:

$$\mathbf{D}_{i,j} := \begin{cases} \deg(v_i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (18.3)$$

Edges

Definition 18.8 Edges

$\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$: Represent some relation between edges and are represented by two-element subset sets of the vertices:

$$e_k = \{v_i, v_j\} \in \mathcal{E} \iff v_i \text{ and } v_j \text{ connected} \quad (18.4)$$

Proposition 18.1 Number of Edges:

A graph G with $n = |\mathcal{V}|$ has between $[0, \frac{1}{2}n(n - 1)]$ edges.

Graph Representations

Adjacency Matrix

Definition 18.9 (unweighted) Adjacency Matrix

\mathbf{A} : Given a graph $G = (\mathcal{V}, \mathcal{E})$ its adjacency matrix is a square matrix $\mathbf{A} \in \mathbb{N}^{n,n}$ defined as:

$$\mathbf{A}_{i,j} := \begin{cases} 1 & \text{if } \exists e(i, j) \\ 0 & \text{otherwise} \end{cases} \quad (18.5)$$

Definition 18.10 weighted Adjacency Matrix

\mathbf{A} : Given a graph $G = (\mathcal{V}, \mathcal{E})$ its weighted adjacency matrix is a square matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ defined as:

$$\mathbf{A}_{i,j} := \begin{cases} \theta_{i,j} & \text{if } \exists e(i, j) \\ 0 & \text{otherwise} \end{cases} \quad (18.6)$$

Diagonal Elements

For a graph without self-loops the diagonal elements of the adjacency are all zero.

Adjacency List

Operations on Graphs

1. Walks

Definition 19.1 Walk: A walk of a graph G as a sequence of vertices with corresponding edges:

$$W = \{v_k, v_{k+1}\}_k^K \in \mathcal{E} \quad (19.1)$$

Definition 19.2 Length of a Walk K : Is the number of edges of that Walk.

2. Paths

Definition 19.3 Path P : Is a walk of a graph G where all visited vertices are distinct (no-repetitions).

Attention: Some use the terms walk for paths and simple paths for paths.

3. Cycles

Definition 19.4 Cycle: Is a path of a graph G where the last visited vertex is the one from which we started.

Types of Graphs

1. Subgraph

Definition 20.1 Subgraph $\mathcal{H} \subseteq \mathcal{G}$: A graph $\mathcal{H} = (U, F)$ is a subgraph of a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ iff:

$$U \subseteq \mathcal{V} \quad \text{and} \quad F \subseteq \mathcal{E} \quad (20.1)$$

2. Components

Definition 20.2 Component: A connected component of a graph G is a connected subgraph of G that is maximal by inclusion – there exist no larger connected containing subgraphs.

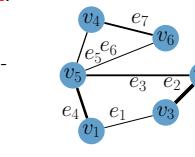
The number of components of a graph G is defined as $c(G)$.

3. Weighted Graph

Definition 20.3 Weighted Graph:

Is a graph G where edges are associated with a weight:

$$\exists \theta_i := \text{weight}(e_i) \quad \forall e_i \in \mathcal{E}$$

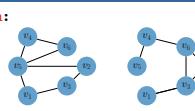


4. Spanning Graph

Definition 20.4 Spanning Graph:

Is a subgraph $\mathcal{H} = (U, F)$ of a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ for which it holds:

$$U = \mathcal{V} \quad \text{and} \quad F \subseteq \mathcal{E} \quad (20.2)$$



5. Connected Graphs

Definition 20.6 (Weakly) Connected Graph:

Is a graph G where there exists a path between any two vertices:

$$\exists P(v_i, \dots, v_j) \quad \forall v_i, v_j \in \mathcal{V} \quad (20.3)$$

Corollary 20.1 Strongly Connected Graph: A directed graph \mathcal{G} is called strongly connected if every node is reachable from every other node.

Corollary 20.2 Components of Connected Graphs: A connected graph \mathcal{G} consist of one component $c(\mathcal{G}) = 1$.

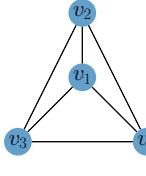
1. Fully Connected/Complete Graph

Definition 20.7 Fully Connected/Complete Graph:

Is a connected graph \mathcal{G} where each node is connected to every other node.

$$\exists e \forall \{v_i, v_j\} \quad \forall v_i, v_j \in \mathcal{V} \quad (20.4)$$

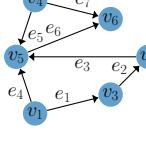
$$|\mathcal{V}| = \frac{1}{2} |\mathcal{V}|(|\mathcal{V}| - 1) \quad (20.5)$$



2. Directed Graphs

Definition 20.8 Directed Graph/Digraph (DG):

A directed graph \mathcal{G} is a graph where edges are direct arcs.



Definition 20.9 Directed Edges/Arcs: Represent some directional relationship between edges and are represented by ordered two-element subset sets of vertices:

$$e_k = \{v_i, v_j\} \in \mathcal{E} \iff v_i \text{ goes to } v_j \quad (20.6)$$

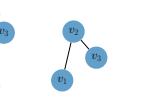
Acyclic Graphs

Definition 21.1 Acyclic Graphs: Are graphs where no cycles exist.

Forests

Definition 21.2 Forests:

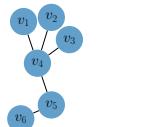
Are acyclic graphs:



Trees

Definition 21.3 Trees:

Are acyclic graphs that are connected.



Binary Trees

Definition 21.4 Binary Tree:

Is a tree where each node $v_i \in \mathcal{V}$ has up to two children:

$$\deg(v_i) \leq 2 \quad \forall v_i \in \mathcal{V} \quad (21.1)$$

Definition 21.5 Binary Search Tree (BST): Is a binary tree, where the left subtree of a node contains only values smaller than the parent and the right subtree contains only values larger than the parent.

Corollary 21.1 Balanced Binary Search Tree: Is a tree that ensures $\mathcal{O}(\log n)$ time for finding or inserting a node. It is a tree where the number of left and right descendants is roughly equal.

Definition 21.6 Complete Binary Trees: A complete binary tree is a tree in which every node of every level of tree has two children, except the last, to the extent that it has to be filled left to right.

Definition 21.7 Fully Binary Tree: Is a tree where every node has either zero or two children.

Definition 21.8 Perfect Binary Tree: Is a complete binary tree where the last level is also filled, a perfect tree of height n needs to have 2^{n-1} nodes.

Binary Max/Min-Heaps

Definition 21.9 Binary Heap:

Is a complete-binary tree where every parent is smaller/larger (min-heap/max-heap) than its children.

Tries/Prefix Trees

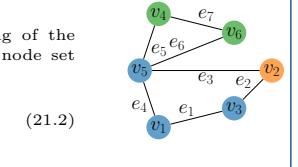
Definition 21.10 Prefix Tree:

Is a tree special kind of tree where each node can have multiple children. It is usually used for prefix lookup of words, where words with the same prefix share the same nodes. It can reduce lookup time from $\mathcal{O}(M \log N)$ for a word of size M with N total words to $\mathcal{O}(M)$. Special terminating nodes are used to indicate if a prefix is an actual word.

1. Graph Layering

Definition 21.11 Graph Layering:

Given a graph G a layering of the graph is a partition of its node set \mathcal{V} into subsets $\{V_1, \dots, V_L\} \subseteq \mathcal{V}$ s.t. $\mathcal{V} = V_1 \cup \dots \cup V_L$ (21.2)



2. Bisection Algorithms

1. Local Approaches

2. Global Approaches

2.2.1. Spectral Decomposition

Definition 21.12 Graph Laplacian (Matrix)

$\mathbf{L}(G)$: Given a graph with n vertices and m edges has a graph laplacian matrix defined as:

$$\mathbf{L} = \mathbf{A} - \mathbf{D} \quad l_{i,j} := \begin{cases} -1 & \text{if } i \neq j \text{ and } e_{i,j} \in \mathcal{E} \\ 0 & \text{if } i \neq j \text{ and } e_{i,j} \notin \mathcal{E} \\ \deg(v_i) & \text{if } i = j \end{cases} \quad (21.3)$$

Corollary 21.2 title:

2.2.2. Inertial Bisection

Proofs

Model Parameter Estimation

Proof 22.1: 6.10:

$$\begin{aligned} p(\mathbf{X}, \mathbf{y}, \theta) &= \frac{p(\theta | \mathbf{X}, \mathbf{y}) p(\mathbf{X}, \mathbf{y})}{p(\mathbf{y} | \mathbf{X}, \theta) p(\mathbf{X}, \theta)} \\ p(\theta | \mathbf{X}, \mathbf{y}) p(\mathbf{X}, \mathbf{y}) &= p(\theta | \mathbf{X}, \mathbf{y}) p(\mathbf{y} | \mathbf{X}) p(\mathbf{X}) \\ p(\mathbf{y} | \mathbf{X}, \theta) p(\mathbf{X}, \theta) &= p(\mathbf{y} | \mathbf{X}, \theta) p(\mathbf{X}, \theta) \\ &= p(\mathbf{y} | \mathbf{X}, \theta) p(\theta | \mathbf{X}) p(\mathbf{X}) \\ \text{eq. (6.6)} &= p(\mathbf{y} | \mathbf{X}, \theta) p(\theta | \mathbf{X}) \\ &= p(\mathbf{y} | \mathbf{X}, \theta) p(\theta) p(\mathbf{X}) \\ \Rightarrow p(\theta | \mathbf{X}, \mathbf{y}) &= \frac{p(\mathbf{y} | \mathbf{X}, \theta) p(\theta) p(\mathbf{X})}{p(\mathbf{y} | \mathbf{X}) p(\mathbf{X})} \end{aligned}$$

Note

This can also be derived by using the normal Bayes rule but additionally condition everything on \mathbf{X} (where the prior is independent on \mathbf{X})

Deep Learning Submodule

Biological Neural Networks

The human nervous system can be broken down into three stages:



1. The **receptors** collect information from the environment – e.g. photons on the retina.
2. The **effectors** generate interactions with the environment – e.g. activate muscles.
3. The **neural network** processes the incoming and outgoing information and controls the flow of information/activation, which is represented by arrows – feedforward and feedback.

1. Neurons/Nerve Cells

Definition 23.1 Neurons/Nerve Cells:

Are the basic building blocks of a neural networks, that processes and transmits information through electrical and chemical signals.

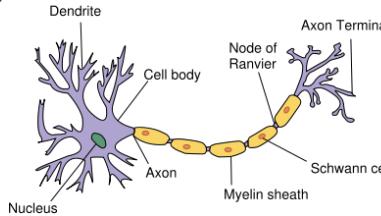


Figure 10: Neuron/Nerve Cell

Definition 23.2 Axon: Signals are sent to other neurons via axons.

Definition 23.3 Synapses: Are the end junctions of neurons, that sit at the end of the Axons. They are complex membranes that are responsible for transmitting signals to other cells.

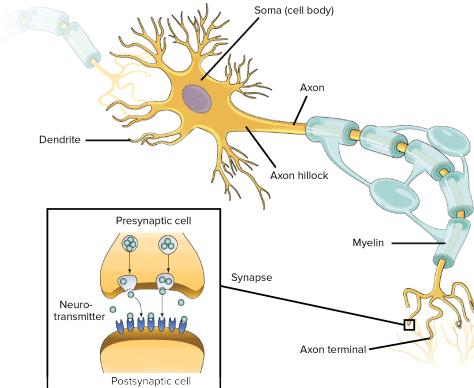


Figure 11: Synapse

Definition 23.4 The Soma and Dendrites:

Are responsible for receiving synaptic signals from other neurons.

Definition 23.5 Action Potential:

Neurons are *electrically excitable*. That means if the net excitation received by a neurons over the dendrites and Soma, is large enough/exceeds a certain *threshold* within a short period of time, then the neuron generates a brief pulse called an *action potential*. In other words the neural cell gets triggered/activated.

Action potentials originate at soma and propagates rapidly along the axon, activating synapses on other neurons as it goes.

Note

Synaptic signals may be excitatory or inhibitory.

2. Artificial Neural Networks (ANN)

Question: how can we model biological neural networks in a mathematical sense?

Idea: start by trying to model a biological neuron.

1. Mathematical Abstraction

We can view a neuron as real valued function that takes n-real valued inputs and maps them to true or false

$$f : \mathbb{R}^n \mapsto \{0, 1\} \quad (23.1)$$

if we take into account that the strength of each input signal may depend on how well the *Axons* transmit then we should include a weight for the Axons:

$$f : \mathbb{R}^n \times \mathbb{R}^d \mapsto \{0, 1\} \quad (23.2)$$

Note

We may also view a neuron as a mapping to a real valued output, if we view the output as firing rate/frequency.

$$f : \mathbb{R}^n \times \mathbb{R}^d \mapsto \mathbb{R} \quad (23.3)$$

History

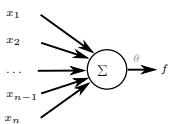
Neurons

1. McCulloch-Pitts (MCP) Neuron

1943

Definition 24.1 McCulloch-Pitts (MCP) Neuron:

Let $x \in \{0, 1\}^n$, $\sigma \in \{-1, +1\}^n$ and $\theta \in \mathbb{Z}$

$$f(x; \sigma, \theta) = \begin{cases} 1 & \text{if } \sum_{i=1}^n \sigma_i x_i \geq \theta \\ 0 & \text{otherwise} \end{cases}$$


Notes

- Is also known as threshold logic unit or linear threshold gate.
- The weights are determined analytically \Rightarrow no learning.
- Has only boolean outputs

1. Logical Units

Any task or phenomenon that can be represented as a logic function can be modeled by a network of MP-neurons.

Proposition 24.1 [example 24.1]

Disjunctive Normal Form of MCPs:

Given: the tuple of the set of all activities $\Sigma = \{\sigma_i\}_{i=1}^n$ and the threshold θ .

Let \mathcal{I} be the set of all subsets I of Σ that activate our neuron:

$$\mathcal{I} = \left\{ I : \sum_{i \in \Sigma} \sigma_i \geq \theta \right\}$$

then the MLP neuron can be written in DNF as:

$$f(x; \sigma, \theta) = \bigvee_{I \in \mathcal{I}} \left(\bigwedge_{i \in I} x_i \wedge \bigwedge_{i \notin I} \neg x_i \right) \quad (24.1)$$

Example 24.1 NAND and DNF:

		$\sum_{i=1}^2 \sigma_i x_i$	f	In order to satisfy NAND we need to set $\sigma_1 = -1$, $\sigma_2 = -1$ and $\theta = -2$.
x_1	x_2			
0	0	0	1	$\sigma_2 = -1$ and $\theta = -2$.
0	1	-1	1	$\mathcal{I} = \{\emptyset, \{1\}, \{2\}\}$
1	0	-1	1	
1	1	-2	0	

$f(x; \sigma, \theta) = (\neg x_1 \wedge \neg x_2) \vee (x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)$

2. Turing Machines

1948

Definition 24.2 Unorganized Machines Alan Turing:

Machines which are largely random in their construction.

Alan Turing asked how to construct unorganized machines, in which organization would emerge from itself.

Definition 24.3 Type A Machines:

Are networks of NAND units, operating in a clocked manner:

	x_1	x_2	y
$y(t+1) = 1 - x_1(t)x_2(t)$	0	0	1
$x_1(t), x_2(t) \in \{0, 1\} \quad \forall t$	1	1	1
	0	0	1
	1	1	0

(24.2)

Note

Alan turing used NAND as all boolean function's can be decomposed into NAND-units.

Definition 24.4 Type B/B Inference (BI) Machines:

Networks, which have modifiers/switches on the connections

Note

Turings original draft was flawed.

3. Willshaw Memory

1968

Definition 24.5 r-Sparse Boolean Vectors:

4. The Perceptron

1958+

Definition 24.6 Perceptron Classifier/Unit: The perceptron is a linear classifier^[def. 4.20] that uses the threshold function in order to classify dichotomies^[def. 4.21]:

$$\hat{y}(\mathbf{x}, \mathbf{w}) = \text{sign}(\mathbf{x}^\top \mathbf{w}) = \begin{cases} +1 & \text{if } \sum_{i=1}^d w_i x_i \geq 0 \\ -1 & \text{else} \end{cases} \quad (24.3)$$

Notes

- The activation function of the perceptron corresponds to the heaviside step function^[def. 6.6, 4.2]
- The perceptron unit is also termed *single-layer-perceptron*, to distinguish it from the misnomer for the multilayer perceptron.

1. Number of Errors

- 4.1.1. Existence of a Separating Solution
- 4.1.2. Rate of Convergence

Definition 24.7 Novikov's Theorem:

- 4.1.3. Uniqueness

5. Learning

1. Hebb's Rule

1949

Units/Activation Functions

Definition 25.1 Activation Pattern: Is the set of neurons of a neural network that are active/inactive for a given input x .

Definition 25.2 Saturation: Corresponds to input value z that causes the maximal activity of the unit (positive or negative).

Proposition 25.1 Zero Gradients Are Problematic:

Zero gradients are often problematic as gradient descent?? updates rely on gradient information thus with zero gradients we do not update any parameters at all.

Threshold Units

1. Linear Units

Has no maximum saturation rate and may hence grow to infinity or minus infinity.

$$\varphi : \mathbb{R}^n \mapsto \mathbb{R} \quad \varphi(\mathbf{x}) = \sum_{i=1}^N w_i x_i + b = \mathbf{w}^\top \mathbf{x} + b$$

$\varphi(\mathbf{x}) = \mathbf{x}$ Activation function is the identity

2. Threshold/Sign Units

Has a maximum saturation rate of 1 and may only take the values one or zero (respective -1, depending on the definition).

$$\varphi : \mathbb{R}^n \mapsto \{0, 1\} \quad \varphi(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^N w_i x_i + b\right) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$$

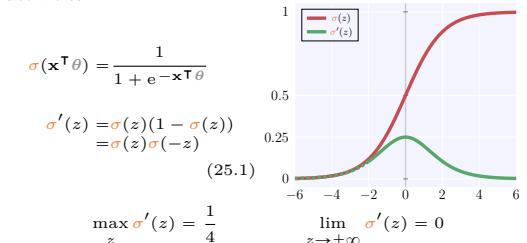
Problem: p.w. continuous constant \Rightarrow zero gradient.

Sigmoid Units

Definition 25.3 Sigmoid Functions: Are functions whose graph is a characteristic "S"-shaped curve.

1. Sigmoid/Logistic Units

Definition 25.4 Sigmoid/Logistic Unit [proof 48] $\sigma : \mathbb{R} \mapsto [0, 1]$: Is a smooth non-linear function that can not have negative activities:



Explanation 25.1 (Sigmoid/Logistic Function).

$$\sigma(\mathbf{x}^\top \theta) = \begin{cases} 0 & -\mathbf{x}^\top \theta \text{ large} \\ 1 & \text{if } \mathbf{x}^\top \theta \text{ large} \\ 0.5 & \mathbf{x}^\top \theta = 0 \end{cases}$$

Pros

- Great for classification problems.

Cons

- Saturated logistic units are non-informative see Corollary 25.4.

Corollary 25.1 Smoothness of σ :
 σ is a smooth function $\sigma \in C^\infty$.

Corollary 25.2 Rotational Symmetry around $(0, 1/2)$:
 $1 - \sigma(z) = \sigma(-z) \iff \sigma(-z) = 1 - \sigma(z)$ (25.2)

Property 25.1 Inverse/Logits [proof 48.2]:
 $\sigma^{-1}(z) = \ln\left(\frac{z}{1-z}\right)$ (25.3)

Notes

- Useful for the output layer in deep neural networks.
- The hill-shaped derivative of the function helps the network when classifying.

Note: σ and pdf

The sigmoid unit corresponds to the CDF of the logistic distribution?? and can be interpreted as a probability. But it is not a real pdf in the mathematical sense. This is because σ approaches 0 on $-\infty$ and 1 on $+\infty$, so that its integral will also be ∞

Logits

Definition 25.5 Logit(s) [proof 48.7]: Is the logarithm of the odds of an event and equals the inverse of the sigmoid function:

$$\text{logit}(p) = \ln\left(\frac{p}{1-p}\right) = \ln(\text{odds}(p)) \quad (25.4)$$

Corollary 25.3 Logits and Neural Networks: The un-normalized and possibly negative inputs to the sigmoid^[def. 25.4]/softmax^[def. 25.7] function are thus called logits.

Vanishing Gradient

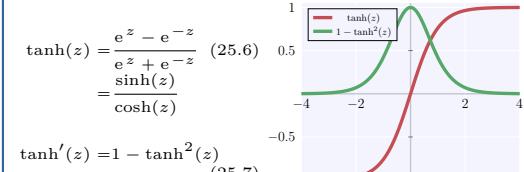
Corollary 25.4 Saturation of logistic units: A saturated logistic unit^[def. 25.4] (large input values z) has a zero/vanishing gradient^[def. 26.11]:

$$\lim_{z \rightarrow \pm\infty} \sigma'(z) = \lim_{z \rightarrow \pm\infty} \sigma(z)\sigma(-z) = 0 \quad (25.5)$$

\Rightarrow parameters will not be updated as the gradient/ e is zero.

2. Hyperbolic Tangent

Definition 25.6 Hyperbolic Tangent/Tanh Unit [proof 48.4] $\tanh : \mathbb{R} \mapsto [-1, 1]$: Is a smooth symmetric non-linear function:



Pros

- Is symmetric around zero with maximal activity at $z \in \{-1, 1\}$
 \Rightarrow i.e. training may lead to pos. or neg. activity (not just pos.)
- Little bit less susceptible to saturation
- Higher derivative values \Rightarrow neuron can better distinguish between similar situations.

Cons

- Still – Saturated tanh units are non-informative.

Property 25.2 Connection to Logistic Unit: $\tanh(z) = 2\sigma(2z) - 1$

Notes

Works almost always better than the sigmoid σ activation function with the one exception for the output layer, if we do binary classification, as the sigmoid function will have a value between 0/1.

3. Swish

Softmax – Generalized Sigmoid Unit

Definition 25.7 Softmax Function [proof 48.5]

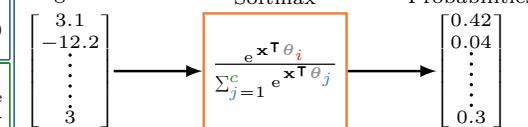
$\sigma^{\max} : \mathbb{R}^c \mapsto \mathbb{R}^c$: Is a non-negative, non-linear, normalized function that generalizes the logistic unit^[def. 25.4] to multiple classes:

$$\sigma_i^{\max}(\mathbf{x}; \Theta) = \frac{\exp(\mathbf{x}^\top \theta_i)}{\sum_{j=1}^c \exp(\mathbf{x}^\top \theta_j)} \quad \Theta = [\theta_1, \dots, \theta_c] \quad (25.8)$$

$$\|\sigma^{\max}\|_1 = \sum_{j=1}^c \sigma_j = 1 \quad \sigma_i > 0 \quad \forall i = 1, \dots, c$$

$$\nabla_{z_j} \sigma_i^{\max}(z) = \begin{cases} \sigma_i^{\max}(z)(1 - \sigma_i^{\max}(z)) & i = j \\ -\sigma_i^{\max}(z)\sigma_j^{\max}(z) & i \neq j \\ = \sigma_i^{\max}(z)(\delta_{ij} - \sigma_j^{\max}(z)) & \end{cases} \quad (25.9)$$

Logits \mathbf{x} Softmax Probabilities



Explanation 25.2.

- The exponential function makes sure that we cannot get negative probabilities.
- If there exist logits^[def. 25.5] of different magnitude, then the logits of small magnitude will basically be ironed out due to the nature of the exponential.
- The name softmax comes from the fact that the function acts as a soft maximum in the sense that the biggest value in magnitude becomes the biggest value while the others pushed towards zero.

Corollary 25.5 Relationship to Sigmoid Unit [proof 48.6]:

$$\text{soft-max with } \theta_1, \theta_2 \iff \text{sigmoid with } \theta := \theta_1 - \theta_2 \quad (25.10)$$

Corollary 25.6 Degrees of Freedom: The problem is over parameterized i.e. more d.o.f. then needed as the probabilities sum to one. Thus we may re-parameterize the problem as:

$$\theta_i \leftarrow \theta_i - \theta_k \quad i = 1, \dots, c \quad \text{s.t.} \quad \theta_k \leftarrow 0 \quad (25.11)$$

2.3.1. Numerical Stable Softmax

Why do we need a stable softmax?

The biggest number a float32 with a 24 bits mantissa can represent is 2^{127} , for an exponent this number is already reach for an input of 88, leading to overflow.

Definition 25.8 Numerical Stable Softmax: Softmax function that avoid over and underflow:

$$\sigma^{\max}(\mathbf{z}) = [\sigma_1^{\max}, \dots, \sigma_c^{\max}] \quad \sigma_k^{\max} \in [0, 1] \quad \forall k$$

$$\sigma_i^{\max}(\mathbf{z}) = \frac{\exp(z_i + a)}{\sum_{j=1}^c \exp(z_j + a)} \quad a = -\max(z_1, \dots, z_c) \quad (25.12)$$

Explanation 25.3. \bullet by subtracting the maximum value from each input, the inputs to the exponents will all be smaller equal to zero, s.t. we cannot have overflow.

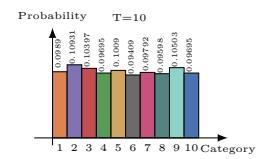
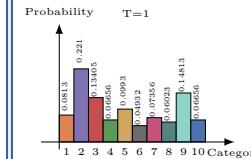
- In addition to that the largest value will be zero such the exponent is one and we have at least one non-zero value in the denominator.

2.3.2. Softmax and Temperature

Definition 25.9 Softmax with Temperature:

It is possible to add a temperature parameter T to the softmax function that allows us to specify how distinct the softmax values are:

$$\sigma_i^{\max}(\mathbf{x}; \Theta) = \frac{\exp\left(\frac{\mathbf{x}^\top \theta_i}{T}\right)}{\sum_{j=1}^c \exp\left(\frac{\mathbf{x}^\top \theta_j}{T}\right)} \quad \Theta = [\theta_1, \dots, \theta_c] \quad (25.13)$$



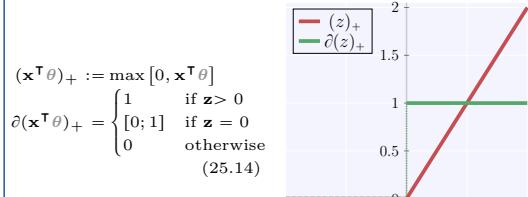
A higher temperature parameter leads to a smoother distribution.

Rectified Units

Definition 25.10 Rectified Units:
Is the use of piece wise linear units C_{pw}^0 .

1. Rectified Linear Units

Definition 25.11 Rectified Linear Units (ReLU) $(\cdot)_+ : \mathbb{R} \mapsto [0; \infty]$: Is the most widely used activation function:



Pros

- Non-linear gradient \Rightarrow easy to optimize.

Cons

- No gradient information for negative entries.
- Dying ReLU problem

Explanation 25.4. ReLU splits the input space into two half spaces:

$$\mathcal{H}_\theta^+ = \{x : \theta^\top x > 0\} \quad \text{and} \quad \mathcal{H}_\theta^- = \{x : \theta^\top x < 0\}$$

separated by the hyperplane \mathcal{H}_θ^0 and constant on $\mathcal{H}_\theta^0 \cup \mathcal{H}_\theta^-$

Corollary 25.7 Dying ReLUs: Refers to the problem that once a ReLU has zero activation it is unlikely to become active again which is due to its gradient:

$$\text{If } z_{lj} = 0 \implies \partial(z_{lj})_+ = 0 \quad (25.15)$$

$$\implies \nabla_{\theta_{lj}} = \nabla_{z_{k-1}} z_{lj} = 0 \quad (25.16)$$

Thus if $z_{lj}(x^t) = 0$ for all inputs of a training sample $X = \{x^1, \dots, x^T\}$ then the ReLU unit:

- is basically dead/the parameters will never be updated and
- will never become active again due to the zero gradient

Corollary 25.8 Activation Pattern:

The activation pattern^[def. 25.1] of a ReLU layer can be described by a heaviside function^[def. 66,42]:

$$H(\Theta x) \in \{0, 1\}^m \quad x \in \mathbb{R}^n \quad (25.17)$$

as each unit will either be active or inactive.

Corollary 25.9 Jacobi Matrix: The Jacobi matrix^[def. 55,6] of a ReLU layer is sparse:

$$\partial F_l := \begin{bmatrix} \tilde{\theta}_{l1}^\top \\ \tilde{\theta}_{l2}^\top \\ \vdots \\ \tilde{\theta}_{l1}^\top \end{bmatrix}, \quad \tilde{\theta}_{lj} = \begin{cases} 0 & \text{if } z_{lj} = 0 \\ \tilde{\theta}_{lj} & \text{else} \end{cases} \quad (25.18)$$

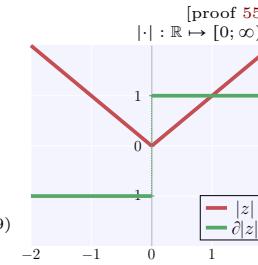
Thus if all units are active we recover a linear network otherwise inactive units are “pruned away”.

2. Absolute Value Units

Definition 25.12 Absolute Value (AbsU)

$$|z| := \begin{cases} z & \text{if } z \geq 0 \\ -z & \text{otherwise} \end{cases}$$

$$\partial(|z|) = \begin{cases} 1 & \text{if } z > 0 \\ [0; 1] & \text{if } z = 0 \\ 0 & \text{otherwise} \end{cases} \quad (25.19)$$



Corollary 25.10 Connection to Absolute value Unit: [proof 48.9]

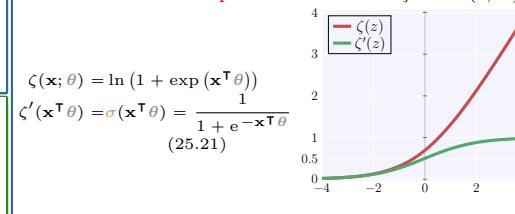
$$(z)_+ = \frac{z + |z|}{2} \iff |z| = 2(z)_+ - z \quad (25.20)$$

Smooth ReLU Approximations

Are no-longer “real” REU units^[def. 25.10] but smooth approximations of the p.w. RELU^[def. 25.11] unit and try to combine the advantages of rectification and smoothness.

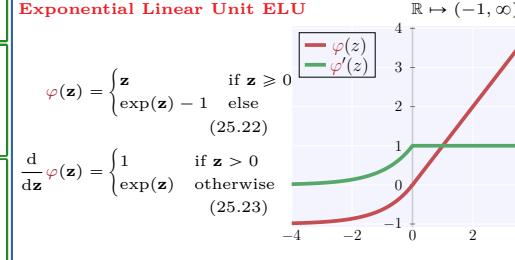
4.0.1. Softplus

Definition 25.13 Softplus



4.0.2. Exponential Linear Unit (ELU)

Definition 25.14 Exponential Linear Unit ELU

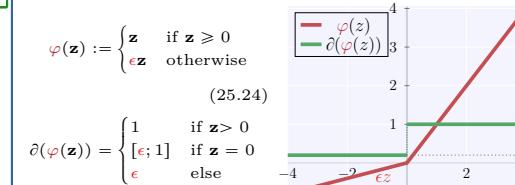


4.0.3. Soft Exponential Linear Unit (SELU)

4.0.4. Parameteric Exponential Linear Unit (PELU)

1. Leaky ReLU (LReLU)

Definition 25.15 Leaky ReLU



2. Approximation Power of pw. Lin. Functions

4.2.1. Hinge Functions

Definition 25.16 Hinge Functions:

Are two hyperplanes^[def. 59,15] that are cut in half and glued together at their intersection:

$$g(x) = \max(\theta_1^\top x + b_1, \theta_2^\top x + b_2) \quad (25.26)$$

Proof 25.2: Property 25.1

$$\sigma\left(\ln\left(\frac{t}{1-t}\right)\right) = \frac{1}{\frac{t}{1-t}} = t \quad (25.29)$$

Proof 25.3: Property 25.2

$$2\sigma(2z) - 1 = 2\frac{1}{1+e^{-2z}} - 1 = \frac{2}{1+e^{-2z}} - \frac{1+e^{-2z}}{1+e^{-2z}} = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \tanh(z)$$

Proof 25.4: Definition 25.6

$$\begin{aligned} \tanh(z) &= \frac{\partial}{\partial z} \sinh(z) = \\ &= \frac{\frac{\partial}{\partial z} \sinh(z) \cdot \cosh(z) - \frac{\partial}{\partial z} \cosh(z) \cdot \sinh(z)}{\cosh^2(z)} \\ &= \frac{\cosh^2(z) - \sinh^2(z)}{\cosh^2(z)} = 1 - \frac{\sinh^2(z)}{\cosh^2(z)} \end{aligned}$$

Proof 25.5: Definition 25.7

Softmax function derivative $\nabla_{z_j} \sigma_i^{\max}(z)$:

$$\begin{aligned} i \neq j : \quad &\nabla_{z_j} \frac{\exp(x^\top \theta_i)}{\sum_{k=1}^c \exp(x^\top \theta_k)} = \\ &Q.R. \quad 0 - e^{x_j} e^{x_i} = -\frac{e^{x_j}}{\sum_{k=1}^c e^{x_k}} \cdot \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} \\ &= -\sigma_i^{\max}(z) \sigma_i^{\max}(z) \\ i = j : \quad &\nabla_{z_i} \frac{\exp(x^\top \theta_i)}{\sum_{k=1}^c \exp(x^\top \theta_k)} = \\ &= \frac{\exp(x^\top \theta_i)}{\left(\sum_{k=1}^c e^{x_k}\right)^2} \\ &= \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} - \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} \cdot \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} \\ &= \sigma_i^{\max}(z) \cdot (\sigma_i^{\max}(z))^2 \end{aligned}$$

Proof 25.6: Corollary 48.6

$$\begin{aligned} \sigma_1^{\max}(x) &= \frac{e^{-x^\top \theta_1}}{e^{x^\top \theta_1} + e^{x^\top \theta_2}} = \frac{e^{-x^\top \theta_1}}{e^{-x^\top \theta_1} + e^{x^\top \theta_2}} \cdot \frac{e^{x^\top \theta_1}}{e^{x^\top \theta_1} + e^{x^\top \theta_2}} \\ &= \frac{1}{e^{x^\top(\theta_1-\theta_2)} + e^{x^\top(\theta_2-\theta_1)}} = \frac{1}{1 + e^{-x^\top \theta}} \end{aligned}$$

Proof 25.7: Definition 25.5

$$\begin{aligned} y &= \frac{e^x}{e^x + 1} = \frac{1}{1 + e^x} = \frac{1 + e^x - 1}{1 + e^x} = 1 - \frac{1}{1 + e^x} \\ &\Rightarrow \frac{1}{1 + e^x} = \frac{1}{1 - y} = \frac{1 - y + y}{1 - y} = 1 + \frac{y}{1 - y} \\ e^x &= \frac{y}{1 - y} \quad \Rightarrow \quad x = \ln \frac{y}{1 - y} \end{aligned}$$

Proof 25.8: Corollary 25.10

$$\text{If } z \geq 0 \implies (z)_+ = 2z - z = z \quad \text{if } z < 0 \quad (25.30)$$

$$\implies (z)_+ = 0 \quad (25.31)$$

Definition 25.4

$$\sigma'(z) = \frac{\partial}{\partial z} (1 + e^{-z})^{-1} = - (1 + e^{-z})^{-2} \frac{\partial}{\partial z} (1 + e^{-z})$$

$$\text{Proof 25.1:} \quad = -\frac{1}{(1 + e^{-z})^2} e^{-z}(-1) = \frac{1}{(1 + e^{-z})^2} (1 - e^{-z} - 1)$$

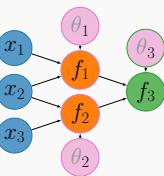
$$= \frac{1}{(1 + e^{-z})} - \frac{1}{(1 + e^{-z})^2} = \sigma(z) - \sigma(z)^2$$

$$\begin{aligned} \sigma'(z) &= \frac{\partial}{\partial z} (1 + e^{-z})^{-1} = - (1 + e^{-z})^{-2} e^{-z}(-1) \\ &= \frac{1}{1 + e^{-z}} \frac{e^{-z}}{1 + e^{-z}} = \frac{1}{1 + e^{-z}} \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}} \\ &= \frac{1}{1 + e^{-z}} \frac{1}{e^z + 1} = \sigma(z) \sigma(-z) \end{aligned}$$

Neural Network Representation

1. DAG Representation

Idea: neural networks that can be represented by acyclic^[def. 74.5] directed^[def. 74.3] graphs^[def. 71.1].



$$\mathbf{x} \mapsto f_3(f_1(x_1, x_2, \theta_1), f_2(x_2, x_3, \theta_2), \theta_3)$$

2. Layered Neural Networks

2.0.1. Neural Network Layers

Definition 26.1 Layered Neural Networks:
Are neural networks where the graph^[def. 71.1] is partitioned into L layers^[def. 26.2].

Definition 26.2 Neural Network Layers $1, \dots, L$:
Are grouped partitions of nodes^[def. 74.8] of the neural network.

1. Composition, Activations and Transfer Maps

Definition 26.3 Transfer Map:

Is a function \mathbf{F}^l that maps the input of one layer $l - 1$ onto the next layer l :

$$\mathbf{F}^l : \mathbb{R}^{n_l} \times \mathbb{R}^{n_{l-1}} \mapsto \mathbb{R}^{n_l} \iff \mathbf{F}^l := \varphi \circ \bar{\mathbf{F}}^l \quad (26.1)$$

with $\bar{\mathbf{F}}^l(\mathbf{x}) = \theta^l \mathbf{x} + b^l \quad b^l \in \mathbb{R}^{m^l}$

Definition 26.4

Composition of Layered Neural Networks:
Is the composition of the L transfer maps^[def. 26.3] $\{\mathbf{F}_l\}_{l=1}^L$ w.r.t. the input, whereas the parameters get concatenated:

$$\mathbf{F} = \mathbf{F}_{L:1} = \mathbf{F}_L \circ \dots \circ \mathbf{F}_1 \quad \mathbf{F}_l : \mathbb{R}^{n_l} \mapsto \mathbb{R}^{n_l+1} \quad (26.2)$$

with $\mathbf{F}(\mathbf{x}; \Theta) = \mathbf{F}_L(\mathbf{F}_{L-1}(\dots \mathbf{F}_1(\mathbf{x}; \theta_1 \dots; \theta_{L-1}; \theta_L)))$
 $\mathbf{z}_l := \mathbf{F}_l \circ \mathbf{F}_{l-1:1}(\mathbf{x}) = \mathbf{F}_l(\mathbf{z}_{l-1})$
 $\mathbf{z}_0 = \mathbf{x} \quad \mathbf{z}_L = \hat{\mathbf{y}}$

2.1.1. Width of a Layer

Definition 26.5 With of a Layer:

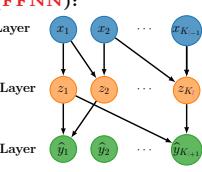
$$\text{width}_l = n_l := \dim(\text{codom}(\mathbf{F}_l)) \quad (26.3)$$

3. Feed Forward Neural Networks

Definition 26.6

Feed Forward Neural Networks (FFNN):

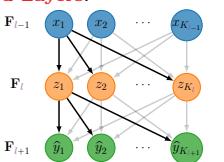
Are neural networks that can be represented by a acyclic^[def. 74.5] directed^[def. 74.3] graphs^[def. 71.1] where information only flows into one direction i.e. along paths^[def. 71.15].



4. Fully Connected Neural Network

Definition 26.7 Fully Connected Layers:

Are neural network layers^[def. 26.2] \mathbf{F}_l whose edges are fully connected^[def. 74.2] to all edges of the previous \mathbf{F}_{l-1} and the next layer \mathbf{F}_{l+1} .



1. Fully Connected FFNN

(FC-FFNN)

Fully Connected Feed Forward Neural Networks :
Are feed forward neural^[def. 26.6] networks that consists only off fully connected layers^[def. 26.7].

4.1.1. Units

Definition 26.9 Units of a Neural Network:
Is the sum of the widths^[def. 26.5] of all layers:

$$\# \text{units} = \sum_{l=1}^L \text{width}_l \quad (26.4)$$

2. Forward Propagation

3. Backpropagation

4.3.1. Derivative of the Network

Definition 26.10 [proof 26.1]

Chain Rule for Neural Networks:

$$\begin{aligned} \partial \mathbf{F} &= \prod_{l=L}^1 \partial \mathbf{F}_l \circ \mathbf{F}_{l-1:1} = \prod_{l=L}^1 \mathbf{J}_{\mathbf{F}_l} \circ \mathbf{F}_{l-1:1} \\ &\equiv \partial \mathbf{F}(\mathbf{x}) = \prod_{l=L}^1 \partial \mathbf{F}_l(\mathbf{z}_{l-1}) \end{aligned} \quad (26.5)$$

with $\mathbf{F}_{0:1} \equiv \text{identity}$ and $\mathbf{z}_{l-1} := \mathbf{F}_{l-1:1}$

4.3.2. Vanishing Gradient Problem

Definition 26.11 Vanishing Gradient Problem: Vanishing Gradient Problem: $\varphi'(\cdot) : \mathbb{R} \mapsto [0, 1]$ hence if we use a deep n-layer network the gradient gets smaller and smaller, as the signal back propagates.

This would mean that the first layer has almost no gradient which would paralyze the network from learning. The reason for this is that we use the chain rule, in order to back propagate the signal and thus multiply all this fractions smaller-equal to one which decreases them exponentially, becoming nearly equivalent to 0.

5. Proofs

Proof 26.1:

$$\begin{aligned} \partial \mathbf{F} &= (\partial \mathbf{F}_l \circ \mathbf{F}_{l-1:1}) \cdot \partial \mathbf{F}_{l-1:1} \\ &= (\partial \mathbf{F}_l \circ \mathbf{F}_{l-1:1}) \cdot (\partial \mathbf{F}_{l-1} \circ \mathbf{F}_{l-2:1}) \cdot \partial \mathbf{F}_{l-2:1} = \dots \end{aligned}$$

Training

Loss Functions

1. Zero-One/Classification Loss
2. Hinge/Perceptron Loss
3. Probabilistic-Log Likelihood Loss

Definition 26.12 (Negative) Log-Likelihood Loss:

Let \mathbf{F} be a model for making probabilistic or noisy predictions $\hat{\mathbf{y}} \in \hat{\mathcal{Y}}$ with $\begin{cases} \hat{\mathbf{y}} \in \mathbb{Z}^m & \text{for probability mass} \\ \hat{\mathbf{y}} \in \mathbb{R}^m & \text{for probability density-functions} \end{cases}$

The likelihood of a label $y_j \in \mathcal{Y}$ given a predicted label \hat{y}_j is given by the negative log-likelihood of the probability distribution $\hat{\mathbf{p}}/f$ of $\hat{\mathbf{y}}$:

$$l_y(\hat{\mathbf{y}}) = -\ln f(y| \hat{\mathbf{y}}; \Theta)$$

Base of Logarithm

The cross entropy uses usually the log 2 and not the natural logarithm but it doesn't really matter as the only vary by a constant factor $\log_2 n = \frac{\ln 2}{\ln n}$.

Categorical-Log Loss

Definition 26.13 [proof 48.10]

Log/Categorical Loss:

Given one-hot^[def. 4.26] encoded hard labels^[def. 4.27] with K individual classes $\{c_1, \dots, c_K\}$.

A categorical output distribution^[def. 66.23] induces a negative log-likelihood loss^[def. 6.4] of:

$$\begin{aligned} l_y(\hat{\mathbf{y}}) &= -\ln f(y|\hat{\mathbf{y}}; \Theta) \\ &= -\delta_{[y=c_k]} \cdot \ln \hat{p}_k \\ &\stackrel{\text{one-hot}}{=} -\sum_{k=1}^K y_k \cdot \ln \hat{p}_k \quad (26.6) \end{aligned}$$

Explanation 26.1. [example 49.1]

Given a true observation/label $y = c_j$ the categorical loss penalizes only the predictive probability $\hat{y}_j \in \hat{\mathbf{p}} = [\hat{y}_1, \dots, \hat{y}_K]^\top$ of the correct c_j by the amount it did not match 1.0.

Cross Entropy Loss

Definition 26.14 [proof 48.12]

Cross Entropy Loss:

Given soft or hard probabilistic targets^[def. 4.28] $\mathbf{p} \in \mathcal{Y}$ with K individual classes $\mathcal{Y} = \{c_1, \dots, c_K\}$ and probabilistic predictions $\hat{\mathbf{p}}$ the cross entropy^[def. 3.7] loss is given by:

$$\text{CE} = l_{\mathbf{p}}(\hat{\mathbf{p}}) = H(\mathbf{p}, \hat{\mathbf{p}}) = -\mathbb{E}_{\mathbf{p}}[\ln \hat{\mathbf{p}}] = H(\mathbf{p}) + \text{KL}(\mathbf{p} \parallel \hat{\mathbf{p}})$$

$$\stackrel{\text{disc.}}{=} -\sum_{k=1}^K p_k \ln \hat{p}_k \quad (26.7)$$

$$\frac{\partial}{\partial z_k} l_{\mathbf{p}}(\hat{\mathbf{p}}(\mathbf{z})) = -\sum_{k=1}^K p_k \frac{1}{\hat{p}(z_k)} \frac{\partial}{\partial z_k} \hat{p}(z_k) \quad (26.8)$$

Attention: $\mathbf{p}/\hat{\mathbf{p}}$ vs. $\mathbf{y}/\hat{\mathbf{y}}$:

This is often ambiguous as if the last layer corresponds to probabilities we can think of:

$$\mathbf{y} = \mathbf{p} \quad \text{and} \quad \hat{\mathbf{y}} = \hat{\mathbf{p}}$$

on the other hand $\hat{\mathbf{y}}$ usually refers to the actual classes:

$$\hat{\mathbf{y}} = \arg \max_{k \in \{1, \dots, K\}} \hat{y}_k \quad (26.9)$$

Attention: Log-loss vs. Cross Entropy Loss:

The cross-entropy loss^[def. 26.14] and the log-loss^[def. 26.13] result in the same loss however the term *log-loss* for the *cross-entropy loss* is a misnomer as the *log/categorical loss* is derived from the *likelihood*, whereas the *cross-entropy loss* is derived from the *cross-entropy*.

In addition to that the *log/categorical loss* only refers to hard labels^[def. 4.27] and does not include soft target^[def. 4.28].

CE Loss with Sigmoid Activation

Corollary 26.1

Cross-entropy/Bernoulli loss with Sigmoid Activation:
If the last layer uses a sigmoid activation^[def. 25.3] the log-loss becomes:

$$l_y(\hat{\mathbf{y}}) = -y \ln \sigma(z) - (1-y) \ln(1-\sigma(z)) \quad (26.10)$$

CE Loss with Softmax Activation

Corollary 26.2 [proof 48.13]

Cross-entropy/Log loss with Softmax Activation:

The cross entropy loss function is usually used together with softmax-activation functions^[def. 25.7] in order to measure how close the output probabilities are to the real distribution. In this case the cross-entropy/log-loss becomes:

$$\begin{aligned} l_y(\hat{\mathbf{y}}) &= -\sum_{k=1}^K y_k \cdot \ln \sigma_k^{\max}(\mathbf{z}) = -\sum_{k=1}^K y_k z_k + \ln \left(\sum_{l=1}^K e^{z_l} \right) \\ &\stackrel{\text{hard targets}}{=} -z_y + \ln \left(\sum_{l=1}^K e^{z_l} \right) \quad (26.11) \end{aligned}$$

$$\frac{\partial}{\partial z_k} l_y(\hat{\mathbf{y}}) = \sigma^{\max}(z_k) - y_k = \hat{p}_k - y_k \quad (26.12)$$

5.3.1. Relationship To Likelihood

Note

Maximizing the likelihood is the same as minimizing the cross entropy.

5.3.2. Binary Cross Entropy (BCE) Loss

Definition 26.15 [proof 48.14]

Binary Cross Entropy (BCE) Loss:

In case of only two classes the cross entropy loss functions can be written as:

$$l_y(\hat{\mathbf{y}}) = -\ln f(y|\hat{\mathbf{y}}; \Theta) = -[y \ln \hat{p}_k + (1-y) \ln(1-\hat{p}_k)] \quad \text{for } y \in \{0, 1\} \quad (26.13)$$

5.3.3. Bernoulli Loss

5.3.4. Squared Loss

Root Mean Squared Error (RMSE)

Mean Absolute Error (MAE)

Mean Squared Error (MSE)

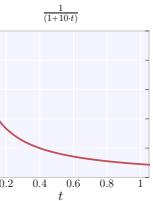
Learning Rates

2. Idea might want to adapt the learning rate w.r.t. the magnitude of the gradient $\mathbf{W}^{(l)}$.

- If the gradient is very/too small \Rightarrow increase learning rate η .
- If the gradient is very/too big \Rightarrow decrease learning rate η .

Time Based Decay

Proposition 26.1 Time Based Decay:



Shrinks the learning rate with time more and more.

$$\eta_t = \eta_0 \cdot \frac{1}{(1 + \text{decay} \cdot t)} \quad (26.14)$$

Proposition 26.2 Time Based Decay after t' iterations:

We may also want to fix the learning rate in the beginning to a small constant C_1 and then slowly decrease it after t' iterations:

$$\eta_t = \min(C_1, C_2) = \min\left(C_1, \frac{C_1 t'}{t}\right) \quad (26.15)$$

Step Based Decay

Proposition 26.3 Step Based Decay:

Divide the learning rate by a fixed factor f after every t' epochs:

$$\eta_t = \eta_0 \cdot f^{\left(\left\lfloor \frac{t}{t'} \right\rfloor\right)} \quad (26.16)$$

Exponential Based Decay

Metrics

BLEU

Trainable Neural Network Layers

Linear Layers

Convolutional Layers

Normalization

Non-trainable Neural Network Layers

Normalization

Layer Normalization

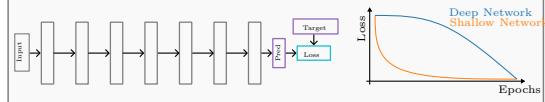
Definition 29.1 Layer Normalization:

Batch Normalization

Residual Layers

Intro

When using neural networks we usually initialize the weights of the layers to some small value. This does not constitute a problem when using shallow networks.



But if we use deep networks initially the information passed through gets multiplied by many random weight matrices s.t. almost none of the output is an actual signal related to the input. Another problem is that if we backpropagate the Jacobin's back to the earlier layers we will multiply many matrices with a small value leading to a vanishing gradient. **Idea:** we can avoid the vanishing gradient if we add the identity to the gradient, which is equal to adding the input to the output:

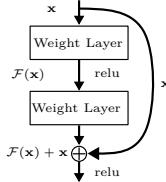
$$\mathcal{F}_{\text{res}}(\mathbf{x}) = \mathcal{F}(\mathbf{x}) + \mathbf{x} \quad \mathbf{J}_{\mathcal{F}_{\text{res}}} = \mathbf{J}_{\mathcal{F}} + \mathbf{I}$$

Definition 29.2 Residual Layers/Blocks

[?]:

A residual layer adds the *input* of a layer to the *output* of another layer:

$$\mathcal{F}_{\text{res}}(\mathbf{x}^{(l)}) = \mathcal{F}(\mathbf{x}^{(l)}) + \mathbf{x}^{(l')} \quad l < l' \quad (29.1)$$



Explanation 29.1 (Why Residual Block). *The residual block does no-longer need to learn everything about the input that needs to be passed along but rather to figure out what information the network can add on-top of the network i.e. we try learn the residual of the input.*

Corollary 29.1 Non-matching dimensions:

If the dimension are not matching we can multiply with an additional weight matrix \mathbf{W} that matches the codomain of \mathcal{F} :

$$\mathcal{F}_{\text{res}}(\mathbf{x}^{(l)}) = \mathcal{F}(\mathbf{x}^{(l)}) + \mathbf{W}\mathbf{x}^{(l')} \quad l < l'$$

Note

Usually multiple layer are skipped by a skip connection, to so that the block is flexible enough to learn the residual.

Why not use concatenation?

If we have more than one or two skip-connections then concatenation may lead to a very large vector.

Note

We need to have a non-linearity, otherwise we add a linearity to a linearity which does not change anything.

Recurrent Neural Networks (RNN)

Goal

Given an observation sequence of inputs $\mathbf{x}^1, \dots, \mathbf{x}^s$ we want to learn the transition^[def. 12.5] output^[def. 12.6]-model of a discrete time state space model^[def. 12.7] that:

(1) Is Markovian^[def. 13.1]

(2) Is Stationary Time-homogeneous/-invariant^[def. 13.15].

$$\mathbf{h}^{(t)} = \mathbf{F}(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \theta_{\mathbf{F}}) \quad \mathbf{h}^0 = \mathbf{0} \quad t = 1, \dots, s \quad (29.2)$$

Vanilla RNNs

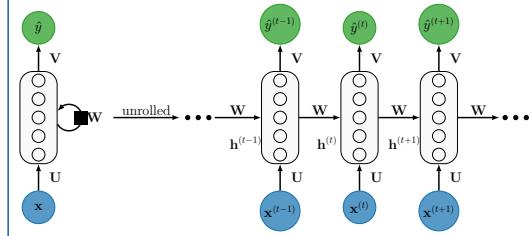
Definition 29.3 Recurrent Neural Network (RNN): Is a linear discrete time state space model^[def. 12.6] with element-wise non-linearities:

State Transition Model:

$$\begin{aligned} \mathbf{F}(\mathbf{h}, \mathbf{x}; \theta_{\mathbf{F}}) &= \mathbf{W}\mathbf{h} + \mathbf{U}\mathbf{x} + \mathbf{b} & \theta_{\mathbf{F}} &= (\mathbf{U}, \mathbf{W}, \mathbf{b}) \\ \mathbf{F} &= \varphi \circ \mathbf{F} \end{aligned} \quad (29.3)$$

Output Model:

$$\begin{aligned} \mathbf{H}(\mathbf{h}; \theta_{\mathbf{H}}) &:= \mathbf{V}\mathbf{h} + \mathbf{c} & \theta_{\mathbf{H}} &= (\mathbf{V}, \mathbf{c}) \\ \mathbf{y} &= \overline{\mathbf{H}}(\mathbf{h}; \theta_{\mathbf{H}}) = \varphi_{\mathbf{H}} \circ \mathbf{H} \end{aligned} \quad (29.4)$$



Main difference to FCFFN^[def. 26.8]

- The weights/parameters of the different layers are *shared!*
- Inputs \mathbf{x}_i are processed in a sequence.

Cons

- Non-parallization:* RNNs cannot parallelized due to the sequential dependency of the stateseq. (29.2).
- Linear-interaction Distance:* it can be difficult for RNNs to relate to related tokens to each other, if there have been a lot of tokens inbetween them.

1. Backpropagation Through Time (BPTT)

Gated Recurrent Units (GRUs)

Long Term Short Term Memory (LSTM)

Types of RNN Tasks

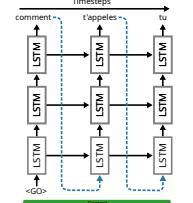
Sequence-to-Sequence (Seq2Seq) Models

Sequence to Sequence Modeling

Sequence-to-sequence modeling is the task of taking an input and producing an output sequence:

- Translation:** taking an input sentence in one-language and outputting it in another.
- Summarization:** taking an input sentence and producing a summary of it.
- Conversation:** taking an input sentence or question and producing an answer to it.

Proposition 32.3 Stacked RNNs:



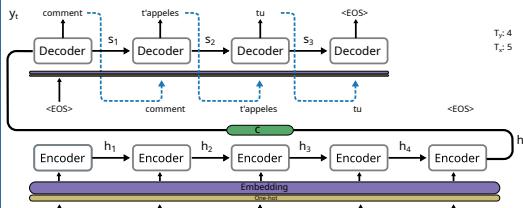
As the input is really difficult to compress and arbitrary length sequence into one context vector of fixed size, the encoder usually consists of stacked RNNs.

Definition 32.1 Seq2Seq Models:

Sequence-to-sequence models consist of encoder-decoder architectures that can be trained in an end-to-end fashion.

- Encoder:** encodes an input sequence $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$ into a context vector \mathbf{c} . The most common approach is an RNN and some nonlinear function q :
$$h_t = f(x_t, h_{t-1}) \quad \mathbf{c} = q(\{h_1, \dots, h_T\}) \quad \text{i.e. } [?] = h_T$$
- Decoder:** uses the context vector to generate an output sequence:

$$\begin{aligned} p(\mathbf{y}) &= p(y_1, \dots, y_T) = \prod_{t=1}^T p(y_t | \{y_1, \dots, y_{t-1}\}, \mathbf{c}) \\ \text{RNN} &= g(y_{t-1}, s_t, \mathbf{c}) \end{aligned}$$



Note

We need:

- A start token for the decoder.
- An end of sentence token for the decoder, so that we know when output is finished.

We can use the same token for this i.e. $<\text{EOS}>$ or different tokens i.e. $<\text{START}> / <\text{EOS}>$.

Considerations for the Encoder

Proposition 32.1 Reverse Input Sequence:

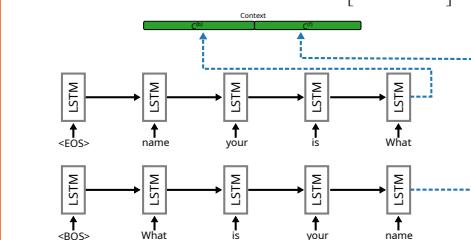
Seq2Seq encoders often process the input sequence in reverse. The benefit is that it is easier for the decoder to focus on the "first" token which is now the last output of the encoder.

Proposition 32.2 Bidirectional RNNs:

Language sentences often have dependencies in both directions, to take care of this one can produce a forward and a backward context and concatenate them before passing them to the decoder:

$$\mathbf{c} = [\mathbf{c}^{(f)} \quad \mathbf{c}^{(b)}] \quad \text{and states} \quad \mathbf{h} = [\mathbf{h}^{(f)} \quad \mathbf{h}^{(b)}] \quad (32.1)$$

$$\mathbf{o} = [\mathbf{o}^{(f)} \quad \mathbf{o}^{(b)}]$$



Autoencoders (AE)

Definition 33.1 Autoencoder

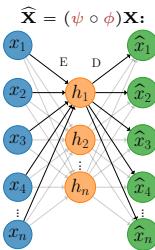
Let $X = \{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^{d \times n}$ and let $\mathbf{X}^\top = [\mathbf{x}_1 \dots \mathbf{x}_n] \in \mathbb{R}^{n \times d}$ be the transposed design matrix.

An auto encoder is defined as:

Encoder $\phi : \mathbb{R}^d \mapsto \mathbb{R}^p$

Decoder $\psi : \mathbb{R}^p \mapsto \mathbb{R}^d$

with $\hat{\mathbf{x}} = (\psi \circ \phi)\mathbf{x}$



$$\begin{aligned}\theta = \{\phi, \psi\} &= \arg \min_{\phi, \psi} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - (\psi \circ \phi)\mathbf{x}_i\|^2 \\ &= \arg \min_{\phi, \psi} \frac{1}{2n} \|\mathbf{X}^\top - (\psi \circ \phi)\mathbf{X}^\top\|_F^2 \quad (33.1)\end{aligned}$$

Explanation 33.1 (Definition 33.1). An autoencoder is a neural network that is supposed to learn the identity map by first encoding and then decoding the input.

1. Linear Autoencoder

Definition 33.2 Variational Autoencoder $\hat{\mathbf{X}} = \mathbf{DCX}$: Let $X = \{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^{d \times n}$ and let $\mathbf{X}^\top = [\mathbf{x}_1 \dots \mathbf{x}_n] \in \mathbb{R}^{n \times d}$ be the transposed design matrix. An auto encoder is defined as:

$$\begin{aligned}\text{Encoder} \quad \mathbf{C} \in \mathbb{R}^{p \times d} : \mathbb{R}^d \mapsto \mathbb{R}^p \quad \text{s.t.} \quad \hat{\mathbf{x}} = \mathbf{DCx} \\ \text{Decoder} \quad \mathbf{D} \in \mathbb{R}^{p \times d} : \mathbb{R}^p \mapsto \mathbb{R}^d \quad (33.2)\end{aligned}$$

$$\begin{aligned}\theta = \{\mathbf{C}, \mathbf{D}\} &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{DCx}_i\|^2 \\ &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \|\mathbf{X}^\top - \mathbf{DCX}^\top\|_F^2 \quad (33.3)\end{aligned}$$

Note

Equation (33.5) is a degenerated problem and has single non-global minima but many minima that are optimal.

Corollary 33.1 Non-unique solution [proof 48.15]: The solution \mathbf{C}^* , \mathbf{D}^* are global degenerative/non-unique solutions.

1. Variational Autoencoder

Definition 33.3 Variational Autoencoder $\hat{\mathbf{X}} = \mathbf{DCX}$: Let $X = \{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^{d \times n}$ and let $\mathbf{X}^\top = [\mathbf{x}_1 \dots \mathbf{x}_n] \in \mathbb{R}^{n \times d}$ be the transposed design matrix. An auto encoder is defined as:

$$\begin{aligned}\text{Encoder} \quad \mathbf{C} \in \mathbb{R}^{p \times d} : \mathbb{R}^d \mapsto \mathbb{R}^p \quad \text{s.t.} \quad \hat{\mathbf{x}} = \mathbf{DCx} \\ \text{Decoder} \quad \mathbf{D} \in \mathbb{R}^{p \times d} : \mathbb{R}^p \mapsto \mathbb{R}^d \quad (33.4)\end{aligned}$$

$$\begin{aligned}\theta = \{\mathbf{C}, \mathbf{D}\} &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{DCx}_i\|^2 \\ &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \|\mathbf{X}^\top - \mathbf{DCX}^\top\|_F^2 \quad (33.5)\end{aligned}$$

Note

Equation (33.5) is a degenerated problem and has single non-global minima but many minima that are optimal.

Corollary 33.2 Non-unique solution [proof 48.15]: The solution \mathbf{C}^* , \mathbf{D}^* are global degenerative/non-unique solutions.

Generative Adversarial Networks (GAN)

Definition 34.1 Gan:

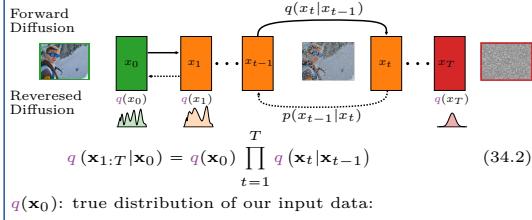
$$\min_G \max_D \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{x \sim p_z(z)} [1 - D(\log D(x))] \quad (34.1)$$

VQ-GAN

Diffusion Models

Definition 34.2 Diffusion Model:

Generative Diffusion Models are models that introduce systematic noise in an iterative process through a Markov Chain and then try to learn to reverse this process in order to generate samples from the underlying distribution:



Definition 34.4 Reverses Diffusion Process:
 $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$ (34.6)

Latent Diffusion Models

History

Definition 34.3 [proof 48.17]

Forward Diffusion Process:

The forward diffusion process incrementally adds noise to the input:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}\left(\sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}\right) \quad \{\beta_t\}_{t=1}^T \in (0, 1) \\ \beta_1 < \beta_2 < \dots < \beta_T \\ \mathbf{x}_t = \sqrt{1-\beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\epsilon \quad \epsilon \sim \mathcal{N}(0, 1) \quad (34.3)$$

The level of added noise is increasing slowly with each time step, regulated by the schedule $\beta_t = \beta_t(t)$ in order to:

- Bring the mean of each new Gaussian closer to zero.
- Limits the rate of variance increase, we want to learn gradually and don't learn anything from pure noise.

$$\lim_{T \rightarrow \infty} q(\mathbf{x}_{1:T}|\mathbf{x}_0) \approx \mathcal{N}(0, \mathbf{I}) \quad (34.4)$$

One Step Forward Process:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}\right) \quad \bar{\alpha}_t := \prod_{s=0}^t \alpha_s \quad (34.5) \\ \mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + (1 - \bar{\alpha}_t)\epsilon \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

Explanation 34.1.

One Step Forward Diffusion Step :

Sampling from a Gaussian and applying eq. (34.3) repeatedly to obtain $q(\mathbf{x}_t|\mathbf{x}_0)$ using eq. (34.2) is expensive, however using a re-parameterization trick we can directly compute $q(\mathbf{x}_t|\mathbf{x}_0)$ without the need to have to apply eq. (34.2).

Notes

If the step-sizes β are too large it becomes difficult to learn the de-noising steps of the reverse process.

Problem

Ideally we would like to calculate $q(\mathbf{x}_{t-1}|\mathbf{x})$ but this is not feasible from section 9 we know that:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}) q(\mathbf{x}_{t-1})}{q(\mathbf{x}_t)} \\ q(\mathbf{x}_t) = \int q(\mathbf{x}_t|\mathbf{x}_{t-1}) q(\mathbf{x}_{t-1}) d\mathbf{x}$$

the integral's to calculate $q(\mathbf{x}_t)$ resp. $q(\mathbf{x}_{t-1})$ are most likely intractable. However if the forward noise step $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ is small, then there is not so much ambiguity about $q(\mathbf{x}_{t-1})$ s.t. we may model $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ by a uni-modal Gaussian distribution.

Idea: replace $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ by a trainable neural network $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$.

Intuition Why This true

For infinitesimal small step-sizes we can convert the forward process into a SDE using Taylor expansion. This SDE can be reverse.

Attention Based Architectures

Definition 34.5 Attention:

Attention is a mechanism that gives a model the knowledge of the relevance of one entity w.r.t. to another entity.



Attention in RNNs

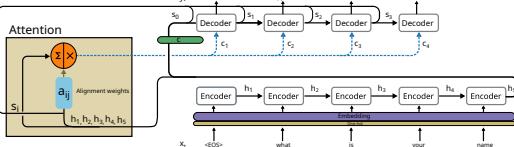
Seq2seq models section 32 have the problem that the decoder takes one context vector as input and all following states depend on this single context vector \mathbf{c} . This does not work well for long sequences.

Idea: create one context vector \mathbf{c}_t for each state $\{\mathbf{s}_t\}_{t=1}^{T_y}$ of the decoder, that describes how relevant / well-aligned each input value $\{\mathbf{x}_i\}_{i=1}^{T_x}$ is with output \mathbf{y}_t .

Definition 34.6 Neural Machine Translation (NMT):

Uses an attention mechanism to add relevant context to each of the output states, as a linear combination of the input states:

$$\begin{aligned}\mathbf{c}_t &= \sum_{t=1}^{T_x} \alpha_{i,t} \mathbf{h}_t = \alpha_{i,t}^T \mathbf{H} \\ &\triangleq \text{alignment}(\mathbf{y}_i, \mathbf{x}_t) \\ \alpha_{i,t} &= \text{softmax}(e_{i,t}) = \frac{\exp(e_{i,t})}{\sum_{j=1}^{T_x} \exp(e_{i,j})} \\ e_{i,k} &= \text{score}(\mathbf{s}_{i-1}, \mathbf{h}_k) = \text{align}(\mathbf{s}_{i-1}, \mathbf{h}_k) \\ \mathbf{s}_i &= f(\mathbf{s}_{i-1}, \mathbf{y}_{i-1}, \mathbf{c}_i)\end{aligned}\quad (34.7)$$



Explanation 34.2.

- We simply take a linear combination of the hidden input states, where each input state is weighted by alignment or score function that describes, how well a input state correlates to an output state.
- The softmax scales the weights s.t. they sum to one and are in the range of [0, 1].

Corollary 34.1 Calculating the output

[?] uses a fully-connected feed forward neural network with a tanh activation function and a softmax layer:

$$\begin{aligned}\tilde{\mathbf{s}}_i &= \tanh(\mathbf{W}_c [\mathbf{c}_i; \mathbf{s}_i]) \\ \mathbb{P}(\mathbf{y}_i | y_i < i, \mathbf{x}) &= \text{softmax}(\mathbf{W}_s \tilde{\mathbf{s}}_i)\end{aligned}\quad (34.8)$$

Corollary 34.2 Initializing s_0 :

Notes

- [?] uses a bidirectional lstm and concatenates the input states as: $\mathbf{h}_i = [\mathbf{h}_i^{(f)}; \mathbf{h}_i^{(b)}] \quad i = 1, \dots, T_x$
- the alignment score function can be used to create a matrix of correlations between input and output tokens.

Alignment-Score Functions

1. Additative Attention

Definition 34.7 Linear Attention:

$$\text{score}(\mathbf{s}_i, \mathbf{h}_j) = \mathbf{v}_a^T \tanh(\mathbf{W}_a \mathbf{s}_i + \mathbf{U}_a \mathbf{h}_j) \quad (34.9)$$

Definition 34.8 Concat Attention:

$$\text{score}(\mathbf{s}_i, \mathbf{h}_j) = \mathbf{v}_a^T \tanh(\mathbf{W}_a [\mathbf{s}_i; \mathbf{h}_j]) \quad (34.10)$$

Trainable Parameters: $\mathbf{W}_a, \mathbf{U}_a, \mathbf{v}_a$

Pros

- Additive attention works better on larger input dimensions.

2. Multiplicative Attention

Definition 34.9 Bilinear Attention:

$$\text{score}(\mathbf{s}_i, \mathbf{h}_j) = \mathbf{s}_i^T \mathbf{W}_a \mathbf{h}_j \quad (34.11)$$

Trainable Parameters: \mathbf{W}_a

Definition 34.10 Dot-Product Attention:

$$\text{score}(\mathbf{s}_i, \mathbf{h}_j) = \mathbf{s}_i^T \mathbf{h}_j \quad (34.12)$$

Trainable Parameters: none

Definition 34.11 Scaled-Dot-Product Attention:

$$\text{score}(\mathbf{s}_i, \mathbf{h}_j) = \frac{\mathbf{s}_i^T \mathbf{h}_j}{\sqrt{T_x}} \quad (34.13)$$

Trainable Parameters: none **Note:** scaling is supposed to help for large input values, where the softmax function may have very small gradients, which makes learning hard:

Explanation 34.3. The softmax function is sensitive to very large inputs, which leads to a vanishing gradient, decreasing the rate of convergence. The dot product of two independent random variables \mathbb{R}^d with mean $\mu = 0$ and variance $\sigma^2 = 1$ has itself a mean of 0 but a variance of d . Thus scaling the dot-product with \sqrt{d} leads to a variance of 1 and helps to stop the input of the softmax from growing to large with d .

Definition 34.12 Location-Based Attention:

Only focuses on the target position:

$$\text{score}(\mathbf{s}_i) = \mathbf{s}_i \quad \alpha_{i,j} = \text{softmax}(\mathbf{W}_a \mathbf{s}_i) \quad (34.14)$$

Trainable Parameters: \mathbf{W}_a

Pros

- Is usually faster and more space efficient using efficient matrix-multiplication techniques.

Hard Attention

Definition 34.13 Hard Attention:

In hard attention one does not calculate a linear combination of input states eq. (34.7) like in soft-attention but picks one input state. This can either be deterministic:

$$\mathbf{c}_i = \arg \max_{\alpha_{i,j}} \{ \mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_{T_y} \} \quad (34.15)$$

or probabilistic:

$$\mathbf{c}_i \sim \mathbb{P}(\{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_{T_y}\} | \alpha_{i,j}) \quad (34.16)$$

Attention:

The arg max function is not differentiable and cannot be used in standard backpropagation. Techniques s.a. Monte Carlos sampling or the re-parameterization trick can be used to circumvent this problem

Local Attention

Definition 34.14 Local Attention/predictive alignment:

At any decoder state \mathbf{s}_t , the networks calculates an aligned position \mathbf{p}_t with a fixed window size D :

$$\mathbf{c}_i = \sum_{t=1}^{T_x} \alpha_{i,t} \mathbf{h}_t \quad \forall \mathbf{h}_j \in [\mathbf{p}_i - D, \mathbf{p}_i + D] \quad (34.17)$$

Corollary 34.3 Gaussian Alignment: One problem is how to choose \mathbf{p}_t without using the non-differentiable arg max. One possibility is a Gaussian kernel:

Pros

- Very long sequences can become really expensive, by using a fixed window size we can keep differentiability and reduce the computational complexity.

Monotonic Alignment

Definition 34.15 Monotonic Alignment:

Aligns \mathbf{p}_i exactly with the target sequence $\mathbf{p}_i = i$ with $i \in \{1, \dots, T_y\}$.

Note

- This may result into input tokens not being considered if the target sequence is shorter than the input sequence and the window size D is small.
- This leads to a symmetric block-diagonal matrix.

Predictive Alignment

Definition 34.16 Predictive Alignment:

Predictive alignment takes into account what input tokens are important to the target sequence i.e. we predict \mathbf{p}_i :

$$\mathbf{p}_t = T_x \cdot \text{sigmoid}(\mathbf{v}_p^T \tanh(\mathbf{W}_p \mathbf{h}_i)) \quad \mathbf{p}_i \in [0, T_x]$$

Note

Leads to block-diagonal matrix with some non-symmetric irregularities.

Corollary 34.4 Gaussian Alignment:

One problem is how to choose \mathbf{p}_t without using the non-differentiable arg max. One possibility is a Gaussian kernel.

Self-Attention

Definition 34.17 Self-Attention:

Self attention relates the input sequence to itself and can adapt any input sequence by replacing the target with the input sequence:

$$\mathbf{e}_{i,j} = \text{score}(\mathbf{h}_i, \mathbf{h}_j) \quad (34.18)$$

Key, Value, Query-Attention

Definition 34.18 Input Embedding:

Let \mathcal{V} be a one-hot encoded^[def. 4.26] vocabulary of size $L := |\mathcal{V}|$ and denote by $\mathbf{w}_j \in \mathcal{V}$ tokens of this vocabulary. We use an embedding $\mathbf{E} \in \mathbb{R}^{d \times |\mathcal{V}|}$ to encode tokens:

$$\mathbf{Ew}_i = \mathbf{x}_i \in \mathbb{R}^d$$

Given a sequence of size T_x we can compute the embedding of the whole sequence as:

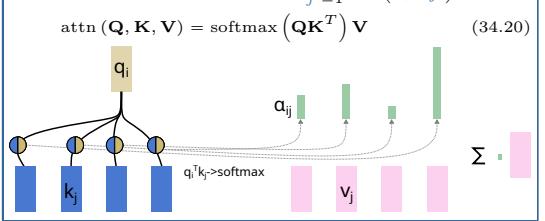
$$\mathbf{w}_{1:T_x} \mathbf{E}^T = \mathbf{X} \in \mathbb{R}^{T_x \times d} \quad \mathbf{w}_{1:T_x} \in \mathbb{R}^{T_x \times |\mathcal{V}|} \quad \mathbf{E} \in \mathbb{R}^{d \times |\mathcal{V}|}$$

Definition 34.19 Query-Key-Value Based Attention: Key-query-value attention, splits the attention mechanism into a key, query and value.

The query \mathbf{q} determines what we are interested in, the key \mathbf{k} what value \mathbf{v} we are currently looking at and the alignment $\alpha_{ij}(\mathbf{q}_i, \mathbf{k}_j)$ determines the weighting of the value:

$$\mathbf{c}_i = \sum_{t=1}^{T_x} \alpha_{i,t} \mathbf{v}_j \quad (34.19)$$

$$\alpha_{i,j} = \text{softmax}(e_{i,j}) = \frac{\exp(\mathbf{q}_i^T \mathbf{k}_j)}{\sum_{j'=1}^{T_x} \exp(\mathbf{q}_i^T \mathbf{k}_{j'})} \quad (34.20)$$



Note: Softmax over matrices

We always take the softmax element-wise and sum over the last dimension.

If \mathbf{A} is a matrix of shape $\mathbb{R}^{l,d}$ we take the softmax as:

$$\text{softmax}(\mathbf{A})_{i,j} = \frac{\exp A_{i,j}}{\sum_{j'=1}^d A_{i,j'}} \quad \forall i \in \{1, \dots, l\}$$

If \mathbf{B} is a tensor of shape $\mathbb{R}^{m,l,d}$ we take the softmax as:

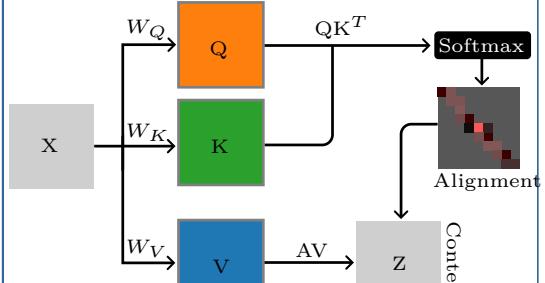
$$\text{softmax}(\mathbf{B})_{k,i,j} = \frac{\exp B_{k,i,j}}{\sum_{j'=1}^d A_{k,i,j'}} \quad \forall i \in \{1, \dots, l\} \quad \forall j \in \{1, \dots, d\}$$

1. Self-Attention

Definition 34.20 Key-Query-Value Self-Attention:

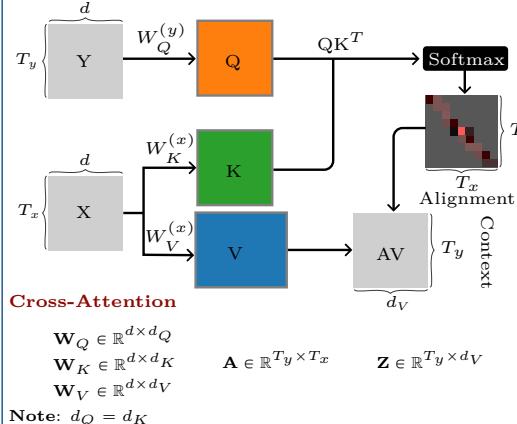
Uses the input sequence for key, value and query:

$$\begin{aligned}\mathbf{q}_i &= \mathbf{W}_q \mathbf{x}_i & \mathbf{k}_i &= \mathbf{W}_k \mathbf{x}_i & \mathbf{v}_i &= \mathbf{W}_v \mathbf{x}_i & \left\{ \begin{array}{l} \mathbf{W} \in \mathbb{R}^{d \times d} \\ \mathbf{X} \in \mathbb{R}^{T_x \times d} \end{array} \right. \\ \mathbf{Q} &= \mathbf{XW}_q & \mathbf{K} &= \mathbf{XW}_k & \mathbf{V} &= \mathbf{XW}_v & \left\{ \begin{array}{l} \mathbf{X} \in \mathbb{R}^{T_x \times d} \end{array} \right.\end{aligned}\quad (34.21)$$



2. Cross-Attention

Definition 34.21 Cross-Attention:
Uses one input sequence $\mathbf{x} \in \mathbb{R}^{T_x \times d}$ for key and value and another $\mathbf{y} \in \mathbb{R}^{T_y \times d}$ for the query, where d corresponds to the embedding dimension:



Explanation 34.4. This helps to relate tokens of the input sequence to the tokens of the output sequence.

Note

d_V may have a different dimension than $d_Q = d_K$ but is usually of the same dimension.

Transformers

[?]

Problems

Transformers remove the RNN component in the encoder-decoder architecture entirely. This comes with solves the previously mentioned problems but also leads to some new problems:

① Positional Relationships:

- The representation of the input tokens/embeddings \mathbf{x} is not position dependent it is simply \mathbf{Ex}_i no-matter what position i .
 - There is also no positional dependence in the self-attention operation, the softmax operation only cares about the word w but not the position i i.e. if we have a word $v \in \mathcal{V}$ at two positions the softmax will be the same.
- ⇒ Need to encode the positional relationships somehow.

② Elementwise Nonlinearity:

- the linear combinationeq. (34.19) is a linear operation. If we stack multiple self-attention layers we end up with a single self-attention layer.
- ⇒ need to add a non-linearity.

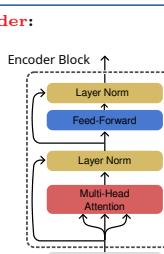
③ Masking:

Given a query \mathbf{q}_i the self-attention mechanism takes into account all tokens $j \{1, \dots, T\}$ but when predicting we only want to consider the past, otherwise predicting becomes trivial.

⇒ we need to mask future entries somehow during the decoding phase.

Definition 34.22 Transformer-Encoder:
The encoder of a transformer encodes the input tokens and adds a positional encoding and then applies a stack of independently parameterized *encoder blocks*, that consist of:

- Multi-head key-query-value self attentioneq. (34.21) block with layer normalization and a residual connection.
 - A feed-forward layer with layer normalization and a residual connection.
- $\mathbf{H}^{(x)} = \text{TransformerEncoder}(\mathbf{X})$



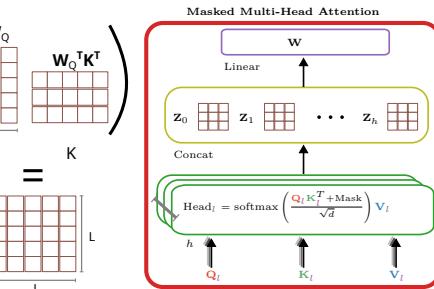
Definition 34.26 Multi-Head Attention:
Multi-head attention is the application of multiple attention heads/mechanisms to one input sequence:

$$\text{MHA}(\mathbf{Q}_l, \mathbf{K}_l, \mathbf{V}_l) = [\text{head}_1, \dots, \text{head}_h] \mathbf{W} \quad (34.22)$$

$$\text{Head}_l = \text{Attention}(\mathbf{X} \mathbf{W}_{q_l}, \mathbf{X} \mathbf{W}_{k_l}, \mathbf{X} \mathbf{W}_{v_l})$$

with $l = \{1, \dots, h\}$.

We can simplify the computation by combining the weight matrices of multiple h attention heads into one weight matrix:
 $\mathbf{W}_{q,k,v} \in \mathbb{R}^{d \times hd}, \quad \mathbf{W}_{q_l,k_l,v_l} \in \mathbb{R}^{d \times d}, \quad l \in \{1, \dots, h\}$



The softmax operates on the same dimension as a normal single head attention operation.

Explanation 34.5.

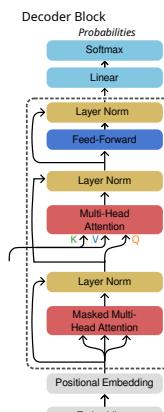
The linear layer after the concatenation has a shared weight-matrix \mathbf{W} which helps in learning a shared representation.

Note

[?] uses $h = 8$ and $\frac{d}{h} = d_Q = d_K = d_V = 64$. This reduces the cost from a single attention mechanism of size $d = 512$ to 64 for the 8 heads.

Efficient Attention

1. Flash Attention



Definition 34.24 The Transformer Block:

Positional Encodings

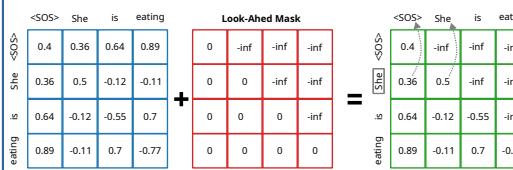
Masking

Definition 34.25 Future Masking:

When we train our encoder we calculate attention scores over the whole sequence. But for the decoder sequence a token y_i should only have access to tokens $k < i$. We can do this by adding a mask to $\mathbf{Q}\mathbf{K}^T$, where $-\inf$ corresponds to entries that are not accessible to a word. After the softmax operation the negative inf value will be zero.

Note: alternatively one can set manually (less efficient) after the softmax:

$$\alpha_{ij} = \begin{cases} \alpha_{ij} & j < i \\ 0 & \text{else} \end{cases}$$



Multi-Head Attention

The attention mechanism learns a linear combination of the valueseq. (34.19). But we can increase the model capacity by using multiple attention mechanisms in parallel.

Graph Networks

1. Undirect Graphs

1. Graph Convolution Networks

1.1.1. General Idea

Given a proximity matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$ of a graph \mathcal{G} and a set of learnable weight matrices $\{\mathbf{W}^l\}_{l=1}^L$ we define the forward propagation as:

$$\mathbf{X}^{l+1} = \varphi(\mathbf{W}^l \mathbf{X}^l) \quad \text{s.t.} \quad \mathbf{X}^L = \mathbf{Y} \quad (36.1)$$

We now want to model the relationship between the different samples a.k.a. nodes by introducing a matrix \mathbf{Q} :

$$\mathbf{X}^{l+1} = \varphi(\mathbf{W}^l \mathbf{X}^l \mathbf{Q}) \quad (36.2)$$

$\mathbf{X}^l \mathbf{Q}$ is a summation over the samples/nodes of \mathbf{X} as $\mathbf{X} \in \mathbb{R}^{d \times n}$.

Question: how should we choose \mathbf{Q} ?

Definition 36.1 Linear Shift Invariant Filter for Graphs: Is a linear function \mathbf{H} over a graph with adjacency matrix \mathbf{A} (possibly degree-normalized $\tilde{\mathbf{A}}$) s.t.

$$\mathbf{H}(\mathbf{A}\mathbf{x}) = \mathbf{A}(\mathbf{H}\mathbf{x}) \iff \mathbf{H} \text{ and } \mathbf{A} \text{ commute} \quad (36.3)$$

Theorem 36.1 Representation of Shift Invariant Filters:

A linear filter \mathbf{H} is \mathbf{A} shift invariant if and only if there exist coefficients $\{\theta_{i=1}\}_{i=1}^n$ s.t. \mathbf{H} can be represented as:

$$\mathbf{H} = \theta_0 \mathbf{I} + \theta_1 \mathbf{A} + \theta_2 \mathbf{A}^2 + \dots + \theta_n \mathbf{A}^n \quad (36.4)$$

1.1.2. Building the Smoothing Matrix

(1) Favor normalized adjacency matrices:

$$\bar{\mathbf{A}} := \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \quad \text{with} \quad \mathbf{D} \text{degree matrix} \quad (36.5)$$

(2) Restrict to linear shift invariant filters of order 1^[def. 36.1]:

$$\mathbf{H}(\theta_0, \theta_1) = \theta_0 \mathbf{I} + \theta_1 \bar{\mathbf{A}} \quad (36.6)$$

(3) Use only one parameter $\theta_0 = \theta_1$

(4) Stack such filters in order to regain high order:

$$\mathbf{H}^k(\theta_0, \theta_1) = \theta^k \mathbf{I} + \theta^k \bar{\mathbf{A}} \quad (36.7)$$

1.1.3. Problems

If we apply an operator and the operator norm^[def. 59.70] is larger than one, then applying the operator repeatedly may lead to blow up.

Lemma 36.1 Eigenvalues of Graph Convolution Layers:

Let $\lambda_1 \geq \lambda_n$ be the eigenvalues of $\bar{\mathbf{A}}$ then it holds that:

$$1 = \lambda_1 \geq \lambda_n \geq -1 \quad \lambda_1, \dots, \lambda_n \in \text{eigenval}(\bar{\mathbf{A}}) \quad (36.8)$$

Corollary 36.1 :

$$\underline{\mathbf{I}} + \bar{\mathbf{A}} \in [0, 2] \quad \Rightarrow \quad \text{possibly blow up} \quad (36.9)$$

1.1.4. Semi-Supervised Classification with GCNNs

Solution to Corollary 1.1.3

Define a degree normalized matrix by using self-loops:

$$\bar{\mathbf{A}} + \mathbf{I} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} + \mathbf{I} \quad \xrightarrow{\text{instead}} \quad \widetilde{\mathbf{D}}^{-\frac{1}{2}} \widetilde{\mathbf{A}} \widetilde{\mathbf{D}}^{-\frac{1}{2}} =: \mathbf{Q}$$

Definition 36.2 Self loop Adjacency Matrix:

$$\widetilde{\mathbf{A}} := \mathbf{A} + \mathbf{I}_n \quad (36.10)$$

Definition 36.3 Self loop Degree Matrix:

$$\widetilde{\mathbf{D}} := \mathbf{D} + \mathbf{I}_n \quad \iff \quad \widetilde{d}_i = 1 + \sum_j^{j < i} \mathbf{A}_{ij} \quad (36.11)$$

Definition 36.4 Coupling Matrix

$\mathbf{Q}:$

$$\mathbf{Q} = \widetilde{\mathbf{D}}^{-\frac{1}{2}} \widetilde{\mathbf{A}} \widetilde{\mathbf{D}}^{-\frac{1}{2}} \quad (36.12)$$

$$q_{ij} = \frac{a_{ij} + \delta_{ij}}{\sqrt{\tilde{d}_i \tilde{d}_j}} \quad \text{wit} \quad \delta_{ij} = \begin{cases} 1 & \text{if } j = i \\ 0 & \text{else} \end{cases}$$

Explanation 36.1 (Coupling Matrix).

$$(\mathbf{X}^l \mathbf{Q})_{ij} = ([\mathbf{x}_1^l \dots \mathbf{x}_n^l] \mathbf{Q})_{ij} = \sum_{k=1}^n x_{ik} q_{kj} \quad (36.13)$$

$$(\mathbf{X}^l \mathbf{Q})_{:j} = \sum_{k=1}^n q_{kj} \cdot \mathbf{x}_k^l \quad (36.14)$$

Thus the coupling matrix^[def. 36.4] leads to a smoothed combination with the neighboring nodes.

Memory Networks

Memories of RNNs consist of hidden states & weights that encode knowledge as dense compressed vectors. However those encodings are usually too small and not compartmentalized enough for large long term storage.

Definition 37.1 Memory Networks: Memory components and a model that is trained to effectively operate with that memory component.

Regularization

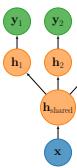
1. Weight Decay

1. L1/Lasso Regularization
2. L2 Regularization
3. Dropout

Multi Task Learning

Definition 39.1 Multi Task Learning:

Share representation of the sub-model across different tasks and learn the model jointly i.e. by a combined objective.



Example 39.1 Invariance to Noise:

- Adding noise to the *inputs*
- Adding noise/small perturbations to the *weights* \Rightarrow regularization of the weights – in variance to local optima.
- Adding noise to the *targets* i.e. probability distribution over the targets (See also section 3 and definition 4.28).

Pretraining/Task Augmentation/Transfer Learning

Definition 40.1 [example 40.1]
Pre-Training: Is the process of pre-training the model on a generic task in order to prime/initialize the model for the actual training.

Example 40.1 Pre-trained Models:

- Pre-trained Word Embeddings
- Pre-trained filters and maps

Embeddings

Positional Encodings

Why do we need a positional encoding?

If we look at the sentences:

Even though she did **not** win the award, she was satisfied.
 Even though she did win the award, she was **not** satisfied.

Then the only differ by the order of one word. In recurrent neural networks the state is recurrent and we feed in the words one after another, but in transformers all words are considered as once and the embedding for a unique word a different positionions has the same embedding i.e. does not contain any positional information.

⇒ We need a way to encode position information for the transformer architecture.

Ways to encode positions

- ① We can use an embedding that encodes the positional relationship.
- ② We can adapt the attention mechanism itself to encode a positional relationship.

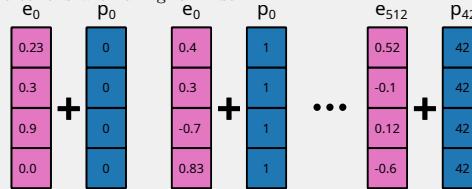
Definition 40.2 Positional Embedding:

Let N be the *maximum possible length of any sequence* and d the input embedding dimension. A positional encoding is a matrix $\mathbf{P} \in \mathbb{R}^{N \times d}$ that encodes the position of a given embedding:

$$\tilde{\mathbf{x}}_i = \mathbf{P}\mathbf{x}_i \quad \mathbf{x}_i = \mathbf{E}\mathbf{w}_i \quad \mathbf{E} \in \mathbb{R}^{d \times |\mathcal{V}|} \quad \mathbf{P} \in \mathbb{R}^{N \times d}$$

Why not just add the token index to the embedding?

If we add the index to the embedding we will have an increasing index that shatters the actual embedding when looking at the tokens with a higher index:



Can we add a normalized index?

If we cannot add the index to the embedding can we add a normalized positional index i.e. $p_j \in [0, 1]$? This does not work either as we cannot use sequences of different length i.e. source and target sequences of different length.

Periodic Embedding

If we cannot use a normalized index, why don't we simply use a periodic embedding with a maximum sequence length?

Sinusoidal (Fourier) Encodings

Definition 40.3

Sinusoidal (Fourier) Positional Encoding:

Idea: we use a period encoding but each dimension of the embedding corresponds to a different wavelength. Let $k = 1, \dots, N$ be the token index and $\delta = 1, \dots, d$ be the embedding index of the corresponding embedding of dimension d . Then the sinusoidal encoding is defined as:

$$PE_k^{(\delta)} = \begin{cases} \sin(\omega_\delta \cdot k) & \text{if } \delta = 2\delta' \\ \cos(\omega_\delta \cdot k) & \text{if } \delta = 2\delta' + 1 \end{cases} \quad \omega_\delta := \frac{1}{10000 \cdot \frac{2\delta}{d}}$$

$$PE_k = \begin{bmatrix} \sin(\omega_1 \cdot k) \\ \cos(\omega_1 \cdot k) \\ \sin(\omega_2 \cdot k) \\ \cos(\omega_2 \cdot k) \\ \vdots \\ \sin(\omega_{d/2} \cdot k) \\ \cos(\omega_{d/2} \cdot k) \end{bmatrix} \in \mathbb{R}^d \quad (40.1)$$

Explanation 41.1. The virtual examples facilitate training by regularization of the original data and thus reduce over fitting.

Note

This may lead to a blow up of the training data but this is no problem due to stochastic gradient descent.

1. Generation By Invariant Transformations

Formula 41.1

[example 39.1]

Invariance Augmentation:

- ① find transformations τ that our original data should be invariant too.
- ② Generate *virtual examples* by applying τ to our original training set:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n \xrightarrow{\tau} \{\tau(\mathbf{x}_i), \mathbf{y}_i\}_{i=1}^n \quad (41.1)$$

Example 41.1 2D Invariances:

- Scale Changes
- Rotations and reflections
- Transformations s.a. PCA

Transfer Learning

Definition 40.4 Transfer Learning: Transfer learning uses learned knowledge of pre-trained models on related tasks or domains.

Knowledge Distillation

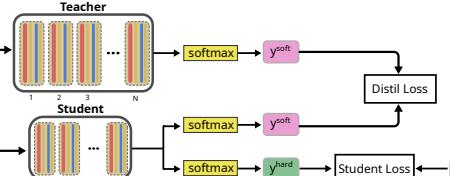
Problem

We have limited hardware resources but small models under-fit to large training datasets. Is there a way to make use of larger pre-trained models?

Definition 40.5

Knowledge Distillation:

Knowledge distillation is the transfer of learned knowledge from a complex teacher model to a simpler student model. The student model is trained to mimic the teacher model by learning from its predictions or soft targets.



The idea is that we train a simpler model not only on the true targets but also try to match it to the distribution of the teacher model predictions by using a *distillation loss*:

$$L = L_{\text{Dist.}}(\mathbf{y}_{\text{Teacher}}, \mathbf{y}_{\text{student}}) + L_{\text{Student}}(\mathbf{y}_{\text{true}}, \mathbf{y}_{\text{student}}) \quad (40.2)$$

Notes

- The distillation loss is usually a loss like the cross entropy loss.
- The softmax values are usually scaled down using a temperature parameter [def. 25.9] T . This is to make it easier to match the student distribution to the teacher distribution. If $T = 1$, then the teacher distribution is relatively non-smooth, almost like a hard target. A softer distribution makes it easier to match our student distribution.

Data Augmentation

Often we have plenty of unlabeled but not so many labeled data. Thus we need a way to create more *artificial/virtual training examples*.

- ① Often we

Definition 41.1 Data Augmentation: Is augmentation of the training data by artificially generated training data.

SOTA Language Models

RNN Models

LSTM-Like Models

Contextualized Word Vectors (CoVe) 08.2017

Intro

While CoVe not an LLM per se CoVe introduced the notion of a pre-trained context that can be used on further downstream tasks. Classical word embedding methods such as Word2Vec and GloVe

Definition 44.1 CoVe [?]

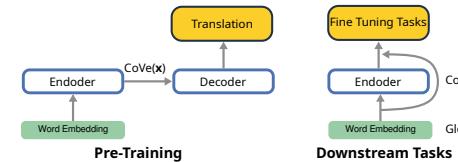
Learned in Translation: Contextualized Word Vectors: CoVe uses an attention-based seq-to-seq NMT^[def. 34.6] model, to create a context vector $\text{CoVe}(\mathbf{x})$ on a pre-training task in order to learn, sentence contexts:

$$\text{CoVe}(\mathbf{x}) = \text{biLSTM}(\text{GloVe}(\mathbf{x})) \quad (44.1)$$

The context vector is then concatenated with an GloVe^{??} embedding to capture the contexts between words:

$$\mathbf{v} = [\text{GloVe}(\mathbf{x}) \quad \text{CoVe}(\mathbf{x})] \quad (44.2)$$

This resulting embedding is than used on downstream tasks s.a. question-awnsing or classification.



Cons

- Pre-training is bound by available supervised training tasks.
- The added benefit of CoVe for the downstream task is limited by the architecture of the downstream task.

ELMO 08.2018

Intro

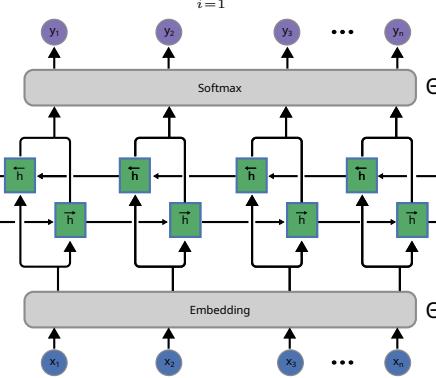
CoVe is limited by the need of supervised training data, while ELMO only needs supervised training data.

Definition 44.2 ELMO [?]

Embeddings from Language Model: ELMO calculates a context aware embedding by calculating a linear combination of the contexts/hidden states \mathbf{h}_i of a stacked bi-LSTM trained on a massive language modeling task:

$$p_{\text{forwd}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n p(\mathbf{x}_i | \mathbf{x}_1, \dots, \mathbf{x}_{i-1}) \quad (44.3)$$

$$p_{\text{backwd}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n p(\mathbf{x}_i | \mathbf{x}_1, \dots, \mathbf{x}_{i-1}) \quad (44.4)$$



The model is trained to minimize the neg. log-likelihood:

$$L = - \sum_{i=1}^n \left(p(\mathbf{x}_i | \mathbf{x}_1, \dots, \mathbf{x}_{i-1}; \Omega_e, \Omega_o, \Omega_{\text{LSTM}}) + p(\mathbf{x}_i | \mathbf{x}_{i+1}, \dots, \mathbf{x}_n; \Omega_e, \Omega_o, \Omega_{\text{LSTM}}) \right) \quad (44.5)$$

Ω_e are the weights of the embedding matrix and Ω_o are the weights of the softmax outputlayer.

The forward and backward hidden-states of one layer are concatenated:

$$\mathbf{h}_k^{(l)} = [h_k^{(l)} \quad \bar{h}_k^{(l)}]$$

s.t. we have $L + 1$ vectors (including) the input embedding per token k :

$$R_k = \{h_k^{(l)} \mid l = 0, \dots, L\}$$

The actual embedding is than a task specific learnable combination:

$$\mathbf{e}_k^{\text{task}} = E(R_k, \Omega^{\text{task}}) = \gamma^{\text{task}} \sum_{j=1}^L s_j^{\text{task}} \mathbf{h}_{k,j} \quad (44.6)$$

s_j^{task} : softmax-normalized learnable weights

Explanation 44.1 (\mathbf{h}_k): The RNN cells have only one state that is shared between different timesteps but we can compute one time-dependent representation of each token of our pre-training language task: \mathbf{R}_k is simply the representation of $w_k \in \mathcal{V}$. If we use the same token multiple times, we can simply use the last representation.

Note

γ^{task} is a scale parameter that is of practical importance to aid the optimization process.

ULMFiT

Transformer Like Models

Encoder-Only Architectures

Encoder only such as BERT learn from left-and right, therefore such architectures are well suited for classification tasks but not as well for generative tasks.

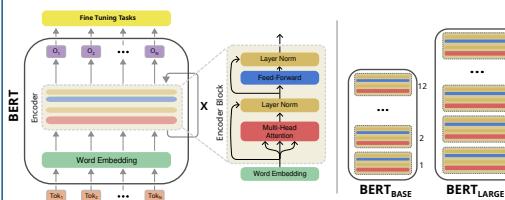
BERT

[110M B/340M L] Param/10.2018

Definition 44.3 [BERT]

Bidirectional Encoder Representations From Transformers:

BERT is a pre-trained LLM that learns contextualized word embeddings and that can be used on many different downstream tasks:



Bert consists on stacked decoder-only transformer blocks. The difference to a normal transformer encoder is the way that BERT is trained.



BERT was trained on Wikipedia (~2.5B words) and Google's BooksCorpus (~800M words). In total ~ 16GB.

Corollary 44.1 Model Size:

BERT comes into pre-trained architectures BERT Base and BERT Large:

	Base	Large
Encoders	12	24
Hidden Size	768	1024
Heads	12	16
Parameters	110M	340M

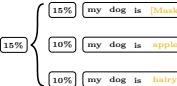
Corollary 44.2 Pre-training BERT:

BERT is pre-trained on two different tasks:

$$L = L_{\text{MLM}} + L_{\text{NSP}}$$

(1) Masked Language Modeling (MLM):

15% of the given words are masked [MASK] and have to be predicted by the model.



This helps BERT to learn a bidirectional representations of the tokens. 10% of the masked words are replaced by a random word and another 10% are replaced by the original word.

This helps for further downstream tasks where we do no longer use the masks.

(2) Next Sentence Prediction (NSP):

BERT is also trained on next sentence prediction, to not only capture relationships between tokens but also between sentences. The sentences are fed into BERT s.t. 50% of the next sentence is correctly following the preceding one.

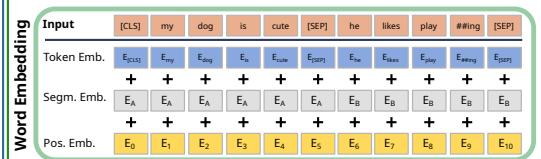
Property 44.1 Training Time:

8xV100x12 days or 280xV100x1day

Corollary 44.3 Word Embeddings:

The word embeddings of BERT consist of three different embeddings, a token embedding, a positional embedding and a segment embedding, that denotes which token belongs to which sentence.

Each sentence starts with a [CLS] token that can be used for downstream classification tasks. Each sentence is separated by a [SEP] token. The length of the input sequence is limited by memory, in the original paper it is limited to 512 tokens.



1. RoBERTa [110M B/340M L] Param/07.2019

Definition 44.4 [ALBERT]:

RoBERTa uses the same architecture as BERT but is only trained on the MLM task, uses a much larger dataset ~ 10x, different training parameters and dynamically changing mask pattern.

RoBERTa is trained on 160GB of data:

- **BOOK CORPUS** (BERT) [5 GB]
- **English Wikipedia dataset** (BERT) [11 GB]
- **CC-NEWS**: 63 million English news articles [2016-2019] [76 GB (after filtering)]
- **OPENWEBTEXT**: Reddit with at least three upvotes [38 GB]
- **STORIES**: Common Crawl data filtered to match the story-like style of Winograd NLP tas [31 GB]

Property 44.2 Training Time:
1024 V100 Tesla GPU for a day.

Pros

- Performs better than BERT in every way.

Cons

- Long training time.

2. DistilBERT 66M Param B/10.2019

Definition 44.5 [DistilBERT]:

Is a distilled version of BERT for fast inference that promises to retain 95% of BERTs performance with 40% fewer parameters and 60% faster training. DistilBERT uses every second encoder block of BERT, with the actual weights of BERT using knowledge distillation^[def. 40.5]. It also uses dynamic masking every other epoch. The distillation loss is given as:

$$L = \frac{L_{\text{problem}} + L_{\text{cross}} + L_{\text{cosine}}}{3} \quad (44.7)$$

$$L_{\text{cross}} = - \sum_{i=1}^n y_i^{\text{(teacher)}} \cdot \log(y_i^{\text{(student)}}) \quad (44.8)$$

$$L_{\text{cosine}} = 1 - \cos(\mathbf{Y}^{\text{(teacher)}}, \mathbf{Y}^{\text{(student)}}) \quad (44.9)$$

Note

The cosine loss is supposed to aligned the directions of the teacher and student.

Property 44.3 Training Time:
8xV100x3.5 days.

3. ALBERT

[B12 /L18/XL60/XX235] M Param

Definition 44.6 [ALBERT]

A lite BERT:

Is a lite version of BERT that has three main differences to BERT:

- ① Factorization of the embedding matrices:

ALBERT reduces the size of the hidden size embeddings.

- ② Cross-layer parameter sharing:

Parameter of different attention and feedforward layers are shared.

- ③ Sentence Order Prediction:

ALBERT replaces the NSP with SOP, which only looks for sentence coherence and not the topic of two sentences that are compared.

ALBERT was trained also only on Wikipedia ($\approx 2.5B$ words) and Google's BooksCorpus ($\approx 800M$ words). In total $\approx 16GB$.

Property 44.4 Training Time:

1024 V100 Tesla GPU for a day.

Pros

- Performs better than BERT in every way.

Cons

- Long training time.

Decoder-Only Architectures

Tend to work well for generative tasks:

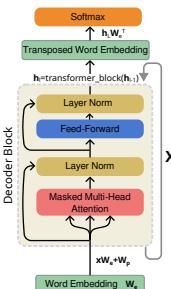
- Question awnsering
- Text generation

GPT-1

117M Param /06.2018

Definition 44.7 [GPT]

Generative Pre-trained Transformer:



Generative Pre-trained Transformer are stacked decoder only transformers. GPT models are autoregressive models that predict one-word at-a-time by minimizing the neg log-likelihood of the next token:

$$L_{LM} = - \sum_{i=1}^k \log p(x_i | x_{i-k}, \dots, x_{i-1})$$

k : context-length

Corollary 44.4 Model Architecture:

Decoders	12
Context Size	512
Hidden Size	768
Batch Size	64

Note

Context Size is the size of the context/window k to be considered for the prediction of the next word.

Corollary 44.5 Training Data:

GPT-1 was trained on:

- BOOK CORPUS (BERT) [5 GB]
- Common Crawl: provides a snapshot of the web and consists of multiple petabytes.

Corollary 44.6 Training Data:

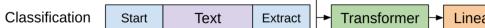
GPT-2 was trained on:

- BOOK CORPUS (BERT) [5 GB]
- Common Crawl: provides a snapshot of the web and consists of multiple petabytes.
- WebText:

12.2019 Supervised Fintuning

The idea of GPT models is to no-longer have a task-specific model but to have a general all purpose model that can be fine-tuned for different tasks.

Lets look at classification, given an input sequence $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ and one label y . For classification one only needs to add linear layer and a softmax activation:



$$P(y|\mathbf{x}_1, \dots, \mathbf{x}_n) = \text{softmax}(\mathbf{h}_L^{(n)} \mathbf{W}_y)$$

The loss to minimize should be a combination of the neg log-likelihood of the true labels and in addition if possible the original LM loss:

$$\begin{aligned} L_{CLS} &= \sum_{(\mathbf{x}, y) \in \mathcal{D}} \log P(y|\mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \sum_{(\mathbf{x}, y) \in \mathcal{D}} \log \text{softmax}(\mathbf{h}_L^{(n)}(\mathbf{x}) \mathbf{W}_y) \\ L_{LM} &= L_{CLS} + \lambda L_{LM} \end{aligned}$$

Optimizing both losses helps in:

- ① Accelerate convergence during training
- ② Improving generalization of the supervised model

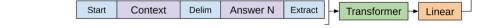
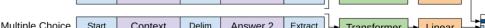
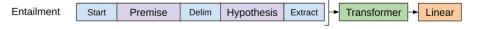
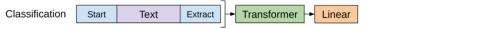


Figure 12: List of different classification problems[?]

Language translation tasks for example can be performed by conditioning on pairs of lang1 = lang2 followed by the target sentence lang1?:

$$P(\text{?} | \text{sen}_1^{(\text{lan1})} = \text{sen}_1^{(\text{lan2})}, \dots, \text{sen}_{\text{target}}^{(\text{lan1})} = \text{?})$$

GPT-2

02.2019

[117 S/345 M/762M L/1542 XL] M Param

Definition 44.8 [GPT-2]

Generative Pre-trained Transformer:

GPT-2 has 10x the size of GPT-1 and comes in four different sizes. The architecture has only some small modifications:

- Layer Normalization is now done before the attention and dense layer.
- Additional Layer Normalization after the final self-attention block.
- Different Initialization

Corollary 44.6 Training Data:

GPT-2 was trained on:

- BOOK CORPUS (BERT) [5 GB]
- Common Crawl: provides a snapshot of the web and consists of multiple petabytes.
- WebText:

Corollary 44.7 Model Architecture:

Decoders	48
Context Size	1024
Hidden Size	1600
Batch Size	512

GPT-3

[/175] B Param /06.2020

Definition 44.9 [GPT-3]

Generative Pre-trained Transformer:

GPT-3 has 10x the size of GPT-2 and uses alternating *dense* and locally *banded, sparse attention* patterns similarly as in *sparse transformers*?. Due to its size, the model is partitioned along both width and depth dimension.

Corollary 44.8 Training Data:

GPT-3 was trained on:

- BOOK CORPUS (BERT) [5 GB]

Common Crawl:

provides a snapshot of the web and consists of multiple petabytes.

Wikimedia:

Books, Articles, and more

Dataset	tokens	Weight in Mix
Common Crawl	410B	60%
WebText2	19B	22%
Books1	12B	8%
Books2	55B	8%
Wikipedia	3B	3%

Corollary 44.9 Model Architecture:

Decoders	96
Context Size	2048
Hidden Size	12288
Batch Size	3.2M

5. XLNet

6. XLM

7. CTRL

BLOOM

Encoder-Decoder Architectures

Work well if we want to create output sequences different from the input sequence:

- Machine Translation
- Text Summarization

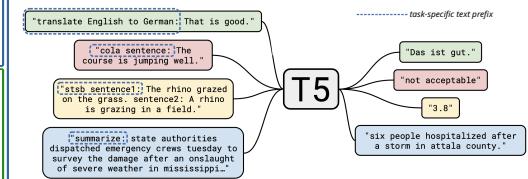
T5

[60M S/220M B/770M L/3B 3B/11B 11B] M Param

Definition 44.10 [T5]

Text-To-Text Transfer Transformer:

Text-To-Text Transfer Transformer does modifies the originally transformer architectures only marginally. T5 rather proposes a unified framework that attempts to combine all language problems into a *text-to-text* format. It uses a short *task-specific text prefix* to distinguish task intentions.



This allows us to use one model, loss and hyperparameters on different downstream tasks. The pre-training tasks consist of sentences with where multiple spawns of tokens are maksed by special tokens:

Orginal : "I", "love", "this", "red", "car".

Input : "I", "<X>", "this", "<Y>"

Target : "<X>", "love", "<Y>", "red", "car", "<EOS>"

The model is fine-tuned for each downstream task separately via "adapter layers" (add an extra layer for training) or "gradual unfreezing". Both fine-tuning approaches only update partial parameters while keeping the majority of the model parameters unchanged.

Explanation 44.2. For regression for example we can train the model to predict text representations of numbers.

Pre-training Tasks

The authors tried various pre-training tasks:

- Language Modeling
- Deshuffling: all words in a sentence are shuffled and the model is trained to predict the original text
- Corrupting Spans: masking a sequence of words from the sentence and the model needs to fill them in.

Corollary 44.10 Training Data:

T5 introduced the C4 — Colossal Clean Crawled Corpus dataset^[def. 46.2].

Unique Transformer Architectures

BART

10.2019

Transformers for Longer Context Lengths

Problem

The input length of vanilla transformers during inference, is upper-bounded by the context length used during training. Due to limited memory and the fact that the context length scales quadratically – both in time $O(L^2d)$ and memory $O(L^2)$ classical transformer architectures have a limited context length between 512 and 1024 tokens.

Predicting in segments of the limited context lengths leads to other problems:

- predicting the first few tokens of a each segment becomes extremely hard
- if we move the segments of size i.e. 1024 by steps of 1, this is very expensive
- long term dependencies can not be captured

1. TransformerXL

186M Param / 06.2019

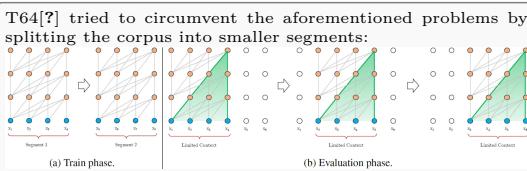


Figure 13: Taken from [?]

This however ignores all contextual information from the previous segments.

TransformerXL tries to solve this problem by introducing a recurrent state and a relative positional encoding. The relative positional encoding is required s.t. that we have positional information of sequences that can be of arbitrary length. The position of the tokens is then encoded as the relative distance between them:

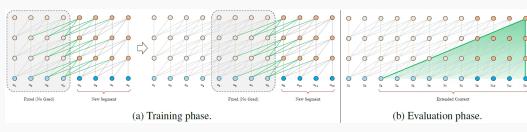


Figure 14: Taken from [?]

Definition 45.1 TransformerXL [proof 48.18]
Transformer Extra Long (CMU/Google AI):

Lets denote two consecutive segments as:

$$\mathbf{s}_\tau = [\mathbf{x}_{\tau,1}, \dots, \mathbf{x}_{\tau,L}] \quad \mathbf{s}_{\tau+1} = [\mathbf{x}_{\tau+1,1}, \dots, \mathbf{x}_{\tau+1,L}]$$

Let n be the n -th hidden layer and define the hidden state of the n -th layer and $\tau + 1$ segment is calculated as:

$$\mathbf{h}_\tau^{n-1} = [\text{SG}(\mathbf{h}_{\tau-1}^{n-1}) \odot \mathbf{h}_\tau^{n-1}] \quad [\odot] \triangleq \text{concat}$$

$$\mathbf{q}_\tau^n = \mathbf{h}_\tau^{n-1} \mathbf{W}_\mathbf{q}^T \quad \mathbf{k}_\tau^n = \tilde{\mathbf{h}}_\tau^{n-1} \mathbf{W}_{\mathbf{k},E}^T \quad \mathbf{v}_\tau^n = \tilde{\mathbf{h}}_\tau^{n-1} \mathbf{W}_\mathbf{v}^T$$

$$\begin{aligned} \mathbf{A}_{\tau,i,j}^n &= \mathbf{q}_{\tau,i}^n \mathbf{k}_{\tau,j}^n && \text{content based addressing} \\ &+ \mathbf{q}_{\tau,i}^n \mathbf{W}_{\mathbf{k},R}^n \mathbf{R}_{i-j} && \text{content dep. positional bias} \\ &+ \mathbf{u}^T \mathbf{k}_{\tau,j}^n && \text{global content bias} \\ &+ \mathbf{v}^T \mathbf{W}_{\mathbf{k},R}^n \mathbf{R}_{i-j} && \text{global positional bias} \end{aligned}$$

$$\mathbf{a}_\tau^n = \text{MaskedSoftmax}(\mathbf{A}_{\tau,i,j}^n) \mathbf{v}_\tau^n$$

$$\mathbf{o}_\tau^n = \text{LayerNorm}(\text{Linear}(\mathbf{a}_\tau^n) + \mathbf{h}_\tau^{n-1})$$

$$\mathbf{h}_\tau^n = \text{Positionwise-Feed-Forward}(\mathbf{o}_\tau^n) \quad \mathbf{h}_\tau^0 := \mathbf{E}_{\mathbf{s}_\tau}$$

SG corresponds to the stop-gradient, meaning that we do not calculate the weights/update the gradient for \mathbf{h}_τ^{n-1} .

$\mathbf{R}_{i-j} \in \mathbb{L}_{\max} \times d$, is a sinusoid encoding matrix with relative distances $i - j$ and no learnable parameters.

Explanation 45.1.

- global bias means that the trainable parameter \mathbf{u} or \mathbf{v} is not dependent on the original query position i see also ?? 48.18.

Pros

- Unlimited context length.
- During inference we can predict as we go and not only segment per segment, which speeds up inference.

3. Compressive Transformer

Training Data

1. Common Crawl

Definition 46.1 [3.1B vewpages/400TiB] 04.2023
Common Crawl (CC): Is a dataset representation of the Internet by removing markup from the HTML pages. CC is extended on a monthly basis by multiple 20TB.

Pros

- Representation of the whole web.

Cons

- Very large.
- A lot of gibberish text like menus, error messages.

1. Colossal Clean Crawled Corpus

Definition 46.2 ≈750GB
Colossal Clean Crawled Corpus (C4): Is a cleaned and filtered version of common crawl from April 2019 it keeps:

- Keep only sentences that end with a punctuation character i.e. period, exclamation mark, question mark, or end quotation mark.
- Remove pages contain explicit words using[?].
- Filtering any lines that contains the word *JavaScript*.
- Pages with placeholder scripts like *lore ipsum*.
- Source code pages are removed by removing all pages that contain curly braces { as many well known programming languages contain those.
- Duplicates are removed by considering three sentence spawns and checking for duplicates.
- [?] is used to only keep enlish pages with a probability of 0.9.

Training Tasks

Proofs

Activation Functions

Definition 25.4

$$\begin{aligned}\sigma'(z) &= \frac{\partial}{\partial z} (1 + e^{-z})^{-1} = - (1 + e^{-z})^{-2} \frac{\partial}{\partial z} (1 + e^{-z}) \\ \text{Proof 48.1: } &= -\frac{1}{(1 + e^{-z})^2} e^{-z} (-1) = \frac{1}{(1 + e^{-z})^2} (1 - e^{-z}) \\ &= \frac{1}{(1 + e^{-z})} - \frac{1}{(1 + e^{-z})^2} = \sigma(z) - \sigma(z)^2 \\ \sigma'(z) &= \frac{\partial}{\partial z} (1 + e^{-z})^{-1} = - (1 + e^{-z})^{-2} e^{-z} (-1) \\ &= \frac{1}{1 + e^{-z}} \frac{e^{-z}}{1 + e^{-z}} = \frac{1}{1 + e^{-z}} \frac{e^{-z}}{e^{-z} + 1} = \sigma(z) \sigma(-z)\end{aligned}$$

Proof 48.2: Property 25.1

$$\sigma\left(\ln\left(\frac{t}{1-t}\right)\right) = \frac{1}{\frac{t}{1-t}} = t \quad (48.1)$$

Proof 48.3: Property 25.2

$$\begin{aligned}2\sigma(2z) - 1 &= 2\frac{1}{1 + e^{-2z}} - 1 = \frac{2}{1 + e^{-2z}} - \frac{1 + e^{-2z}}{1 + e^{-2z}} \\ &= \frac{1 - e^{-2z}}{1 + e^{-2z}} = \frac{e^{-z} - e^{-z}}{e^{-z} + e^{-z}} = \tanh(z)\end{aligned}$$

Proof 48.4: Definition 25.6

$$\begin{aligned}\tanh(z) &= \frac{\partial}{\partial z} \sinh(z) = \\ &= \frac{\frac{\partial}{\partial z} \sinh(z) \cdot \cosh(z) - \frac{\partial}{\partial z} \cosh(z) \cdot \sinh(z)}{\cosh^2(z)} \\ &= \frac{\cosh^2(z) - \sinh^2(z)}{\cosh^2(z)} = 1 - \frac{\sinh^2(z)}{\cosh^2(z)}\end{aligned}$$

Proof 48.5: Definition 25.7

Softmax function derivative $\nabla_{z_j} \sigma_i^{\max}(z)$:

$$\begin{aligned}i \neq j: \quad &\nabla_{z_j} \sum_{k=1}^C \exp(x^\top \theta_k) \\ &\stackrel{\text{Q.R.}}{=} \frac{0 - e^{x_j} e^{x_i}}{\left(\sum_{k=1}^C e^{x_k}\right)^2} = -\frac{e^{x_j}}{\sum_{k=1}^C e^{x_k}} \cdot \frac{e^{x_i}}{\sum_{k=1}^C e^{x_k}} \\ &= -\sigma_i^{\max}(z) \sigma_j^{\max}(z) \\ i = j: \quad &\nabla_{z_i} \sum_{k=1}^C \exp(x^\top \theta_k) \\ &= \frac{e^{x_i} \sum_{k=1}^C e^{x_k} - e^{x_i} e^{x_i}}{\left(\sum_{k=1}^C e^{x_k}\right)^2} \\ &= \frac{e^{x_i}}{\sum_{k=1}^C e^{x_k}} - \frac{e^{x_i}}{\sum_{k=1}^C e^{x_k}} \cdot \frac{e^{x_i}}{\sum_{k=1}^C e^{x_k}} \\ &= \sigma_i^{\max}(z) - (\sigma_i^{\max}(z))^2\end{aligned}$$

Proof 48.6: Corollary 48.6

$$\begin{aligned}\sigma_1^{\max}(x) &= \frac{e^{x^\top \theta_1}}{e^{x^\top \theta_1} + e^{x^\top \theta_2}} = \frac{e^{-x^\top \theta_1}}{e^{-x^\top \theta_1} e^{x^\top \theta_1} + e^{x^\top \theta_2}} \\ &= \frac{1}{e^{x^\top (\theta_1 - \theta_1)} + e^{x^\top (\theta_2 - \theta_1)}} = \frac{1}{1 + e^{-x^\top \theta}}\end{aligned}$$

Proof 48.7: Definition 25.5

$$\begin{aligned}y &= \frac{e^x}{e^x + 1} \cdot \frac{1}{1 + e^{-x}} = \frac{1 + e^x - 1}{1 + e^x} = 1 - \frac{1}{1 + e^x} \\ &\Rightarrow \frac{1}{1 + e^x} = \frac{1}{1 - y} = \frac{1 - y + y}{1 - y} = 1 + \frac{y}{1 - y} \\ e^x &= \frac{y}{1 - y} \quad \Rightarrow \quad x = \ln \frac{y}{1 - y}\end{aligned}$$

Proof 48.14 Binary Cross Entropy Loss^[def. 26.15]:
The binary cross entropy can be derived directly from the Bernoulli Distribution??, where y signifies that we penalize only the true label:

$$\begin{aligned}L_y &= \hat{p}^y \cdot (1 - \hat{p})^{1-y} \quad \text{for } y \in \{0, 1\} \\ l_y &= -\log L_y = -[y \ln \hat{p} + (1 - y) \ln(1 - \hat{p})]\end{aligned}$$

Models

Autoencoders

Proof 48.15 [Corollary 33.2]: Let C^* , D^* be the global solution of problem Equation (33.5). Then we may define a matrix A s.t.:

$$D^* C^* = D^* (A^{-1} A) C^* = \tilde{D}^* \tilde{C}^* \quad (48.6)$$

Proof 48.16 ??:

$$\begin{aligned}DCX &\stackrel{\text{eq. (59.129)}}{=} DCU \Sigma V^H = \underbrace{U_k^T}_{\text{Equation (59.37)}} \underbrace{U_k}_U \underbrace{\Sigma}_{\text{V}} V^H \\ &= U_k \Sigma_m V_m^H = X_m\end{aligned}$$

Diffusion Models

Proof 48.10: Definition 26.13

$$\begin{aligned}l_y(\hat{y}) &= -\log f(y|\hat{y}; \Theta) = -\log \prod_{i=1}^k \delta_{[\hat{y}_i=c_i]} \\ &= -\sum_{i=1}^k \delta_{[\hat{y}_i=c_i]} \cdot \log \hat{p}_i = \delta_{[\hat{y}_i=c_i]} \cdot \log \hat{p}_i\end{aligned}$$

Proof 48.17 One Step Forward Diffusion Model^[def. 34.3]: Let $\{\epsilon\}_{i=1}^T \sim \mathcal{N}(0, 1I)$:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \epsilon_t = \mathbf{x}_{t-1} \sqrt{\alpha_t} + \sqrt{1 - \alpha_t} \epsilon_t \\ &= (\mathbf{x}_{t-2} \sqrt{\alpha_{t-1}} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-1}) \sqrt{\alpha_t} + \sqrt{1 - \alpha_t} \epsilon_t \\ &= \mathbf{x}_{t-2} \sqrt{\alpha_t \alpha_{t-1}} + \sqrt{\alpha_t (1 - \alpha_{t-1})} \epsilon_{t-1} + \sqrt{1 - \alpha_t} \epsilon_t \\ &:= Y \\ Y &\sim \mathcal{N}(0, \alpha_t(1 - \alpha_{t-1})) \quad Z \sim \mathcal{N}(0, 1 - \alpha_t) \\ Y + Z &\stackrel{??}{=} \mathcal{N}(0, \alpha_t(1 - \alpha_{t-1}) + (1 - \alpha_t)) \\ \mathcal{N}(0, 1 - \alpha_t \alpha_{t-1}) &= \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-1} \\ \mathbf{x}_t &= \mathbf{x}_{t-2} \sqrt{\alpha_t \alpha_{t-1}} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-1} \\ &\vdots \\ \mathbf{x}_t &= \mathbf{x}_{t-2} \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_0} + \sqrt{1 - \alpha_t \alpha_{t-1} \cdots \alpha_0} \epsilon_0 \\ \mathbf{x}_t &= \mathbf{x}_{t-2} \sqrt{\alpha_t} + \sqrt{1 - \alpha_t} \epsilon_0\end{aligned}$$

Transformers

Proof 48.18 Relative Positional Encoding^[def. 45.1]:
Relative Positional Encoding:

We can decompose the Vanilla query-key interaction:

$$\begin{aligned}a_{ij} &= q_i k_j^T = (x_i + p_i) W^q ((x_j + p_j) W^k)^T \\ &= x_i W^q W^k x_j^T + x_i W^q W^k p_j^T \\ &\quad + p_i W^q W^k x_j^T + p_i W^q W^k p_j^T\end{aligned}$$

Re-parameterization:

$$p_j \xrightarrow{\text{replace with relative encd.}} r_{i-j} \in \mathbb{R}^d$$

W^k Split into $\begin{cases} W_E^k : \text{Content Information} \\ W_R^k : \text{Positional Information} \end{cases}$

$p_i W^k$ replaced with trainable parameters $\begin{cases} u : \text{Content Information} \\ v : \text{Positional Information} \end{cases}$

This leads to the following entries of the attention matrix:

$$\begin{aligned}a_{ij} &= \underbrace{x_i W^q W_E^k x_j^T}_{\text{content-based addressing}} + \underbrace{x_i W^q W_R^k r_{i-j}}_{\text{content-dep. positional bias}} \\ &\quad + \underbrace{u W_E^k x_j^T}_{\text{global content bias}} + \underbrace{v W_R^k r_{i-j}^T}_{\text{global positional bias}}\end{aligned}$$

Proof 48.13 Deriv. of CE with Softmax^[cor. 26.2] and ?? 48.11:

$$\begin{aligned}\frac{\partial}{\partial z_k} l_y(z) &\stackrel{\text{eq. (26.8)}}{=} -\sum_{j=1}^K y_j \frac{1}{\sigma^{\max}(z_j)} \frac{\partial}{\partial z_k} \sigma^{\max}(z_j) \\ &\stackrel{\text{eq. (25.9)}}{=} -y_k \frac{\sigma^{\max}(z_k)(1 - \sigma^{\max}(z_k))}{\sigma^{\max}(z_k)} \quad j = k \\ &\quad - \sum_{j \neq k} y_j \frac{-\sigma^{\max}(z_k) \sigma^{\max}(z_j)}{\sigma^{\max}(z_j)} \quad j \neq k \\ &= -y_k (1 - \sigma^{\max}(z_k)) + \sum_{j \neq k} y_j \sigma^{\max}(z_k) \\ &= -y_k + \sigma^{\max}(z_k) y_k + \sum_{j \neq k} y_j \sigma^{\max}(z_k) \\ &= \sigma^{\max}(z_k) - y_k = \hat{p}_k - y_k\end{aligned}$$

Examples

Losses

Example 49.1 Categorical Loss^[def. 26.13]:

Given an observation $\mathbf{y} = [0 \ 1 \ 0 \ 0 \ 0]^\top$ the *categorical loss*^[def. 26.13] of the predictive distributions: $\hat{\mathbf{y}} = [\alpha_1, 0.24, \alpha_3, \alpha_4, \alpha_5]$ with $\sum_{i=1}^5 \hat{y}_i = 1$ is equal for any $\alpha_{1,3,4,5} \in [0, 1]$.

Math Submodule

Set Theory

Definition 50.1 Set $A = \{1, 3, 2\}$: is a well-defined group of distinct items that are considered as an object in its own right. The arrangement/order of the objects does not matter but each member of the set must be unique.

Definition 50.2 Empty Set $\{\} / \emptyset$: is the unique set having no elements/cardinality^[def. 50.5] zero.

Definition 50.3 Multiset/Bag: Is a set-like object in which multiplicity^[def. 50.4] matters, that is we can have multiple elements of the same type.
I.e. $\{1, 1, 2, 3\} \neq \{1, 2, 3\}$

Definition 50.4 Multiplicity: The multiplicity n_a of a member a of a multiset^[def. 50.3] S is the number of times it appears in that set.

Definition 50.5 Cardinality $|S|$: Is the number of elements that are contained in a set.

Definition 50.6 The Power Set $\mathcal{P}(S)/2^S$: The power set of any set S is the set of all subsets of S , including the empty set and S itself. The cardinality of the power set is 2^S is equal to $2^{|S|}$.

1. Closure

Definition 50.7 Closure: A set is *closed* under an operation Ω if performance of that operations onto members of the set always produces a member of that set.

2. Open vs. Closed Sets

Definition 50.8 Open Sets:

- Euclidean Spaces:** A subset $U \in \mathbb{R}$ is open, if for every $x \in U$ it exists $\epsilon(x) \in \mathbb{R}+$ s.t. a point $y \in \mathbb{R}$ belongs to U if:

$$\|x - y\|_2 < \epsilon(x) \quad (50.1)$$

- Metric Spaces**^[def. 59.65]: a Subset U of a metric space (M, d) is open if:

$$\exists \epsilon > 0 : \text{if } d(x, y) < \epsilon \quad \forall y \in M, \forall x \in U \implies y \in U \quad (50.2)$$

- Topological Spaces**^[def. 61.2]: Let (X, τ) be a topological space. A set A is said to be open if it is contained in τ .

Definition 50.9 Closed Set: Is the complement of an open set^[def. 50.8].

Definition 50.10 Bounded Set: A set $S \subset \mathbb{R}^n$ is *bounded* if there exists a constant K s.t. the absolute value of every component of every element of S is less or equal to K .

3. Number Sets

1. The Real Numbers

3.1.1. Intervals

Definition 50.11 Closed Interval $[a, b]$: The closed interval of a and b is the set of all real numbers that are within a and b , including a and b :

$$[a, b] = \{x \in \mathbb{R} | a \leq x \leq b\} \quad (50.3)$$

Definition 50.12 Open Interval (a, b) : The open interval of a and b is the set of all real numbers that are within a and b :

$$(a, b) = \{x \in \mathbb{R} | a < x < b\} \quad (50.4)$$

2. The Rational Numbers

Example 50.1 Power Set/Cardinality of $S = \{x, y, z\}$: The subsets of S are:
 $\{\emptyset\}, \{x\}, \{y\}, \{z\}, \{x, y\}, \{x, z\}, \{y, z\}, \{x, y, z\}$ and hence the power set of S is $\mathcal{P}(S) = \{\{\emptyset\}, \{x\}, \{y\}, \{z\}, \{x, y\}, \{x, z\}, \{y, z\}, \{x, y, z\}\}$ with a cardinality of $|S| = 2^3 = 8$.

4. Set Functions

1. Submodular Set Functions

Definition 50.13 Submodular Set Functions: A submodular function $f : 2^\Omega \mapsto \mathbb{R}$ is a function that satisfies:

$$f(A \cup \{x\}) - f(A) \geq f(B \cup \{x\}) - f(B) \quad \forall A \subseteq B \subseteq \Omega \quad \{x\} \in \Omega \setminus B \quad (50.5)$$

Explanation 50.1 (Definition 50.13). Adding an element x to the smaller subset A yields at least as much information/value gain as adding it to the larger subset B .

Definition 50.14 Montone Submodular Function: A monotone submodular function is a submodular function^[def. 50.13] that satisfies:

$$f(A) \leq f(B) \quad \forall A \subseteq B \subseteq \Omega \quad (50.6)$$

Explanation 50.2 (Definition 50.14). Adding more elements to a set will always increase the information/value gain.

2. Complex Numbers

Definition 50.15 Complex Conjugate \bar{z} : The complex conjugate of a complex number $z = x + iy$ is defined as:

$$\bar{z} = x - iy \quad (50.7)$$

Corollary 50.1 Complex Conjugate Of a Real Number: The complex conjugate of a real number $x \in \mathbb{R}$ is x :

$$\bar{x} = x \implies x \in \mathbb{R} \quad (50.8)$$

Formula 50.1 Euler's Formula: $e^{\pm ix} = \cos x \pm i \sin x \quad (50.9)$

Formula 50.2 Euler's Identity: $e^{\pm i} = -1 \quad (50.10)$

Note

$$e^n = 1 \Leftrightarrow n = i 2\pi k, \quad k \in \mathbb{N} \quad (50.11)$$

Sequences&Series

Definition 51.1 Index Set: Is a set^[def. 50.1] A , whose members are labels to another set S . In other words its members index member of another set. An index set is build by enumerating the members of S using a function f s.t.

$$f : A \rightarrow S \quad A \in \mathbb{N} \quad (51.1)$$

Definition 51.2 Sequence $(a_n)_{n \in A}$: A sequence is an an by an index set A enumerated multiset^[def. 50.3] (repetitions are allowed) of objects in which *order* does matter.

Definition 51.3 Series: is an infinite ordered set of terms combined together by addition.

1. Types of Sequences

R 1. Arithmetic Sequence

Definition 51.4 Arithmetic Sequence: Is a sequence where the *difference* between two consecutive terms constant i.e. $(2, 4, 6, 8, 10, 12, \dots)$.

$$t_n = t_0 + nd \quad d : \text{difference between two terms} \quad (51.2)$$

2. Geometric Sequence

Definition 51.5 Geometric Sequence: Is a sequence where the *ratio* between two consecutive terms constant i.e. $(2, 4, 8, 16, 32, \dots)$.

$$t_n = t_0 \cdot r^n \quad r : \text{ratio between two terms} \quad (51.3)$$

Property 51.1 Sum of Geometric Sequence:

$$\sum_{k=1}^n ar^{k-1} = \frac{a(1 - r^n)}{1 - r} \quad (51.4)$$

2. Converging Sequences

1. Pointwise Convergence

Definition 51.6 $\lim_{n \rightarrow \infty} f_n = f$ pointwise
Pointwise Convergence^[?]: Let (f_n) be a sequence of functions with the same domain^[def. 54.8] and codomain^[def. 54.9]. The sequence is said to converge pointwise to its *pointwise limit function* f if it satisfies:

$$|\lim_{n \rightarrow \infty} f_n(x) - f(x)| = 0 \quad \forall x \in \text{dom}(f_i) \quad (51.5)$$

2. Uniform Convergence

Definition 51.7 $\lim_{n \rightarrow \infty} f_n = f$ uniform/ $f_n \xrightarrow{\infty} f$
Uniform Convergence^[?]: Let (g_n) be a sequence of functions with the same domain^[def. 54.8] and codomain^[def. 54.9]. The sequence is said to converge uniformly to its *pointwise limit function* f if it satisfies:

$$\exists \epsilon > 0 : \exists n \geq 1 \sup_{x \in \text{dom}(f_i)} |g_n(x) - f(x)| < \epsilon \quad \forall x \in \text{dom}(f_i) \quad (51.6)$$

Note

Uniform convergence is characterized by the uniform norm^{??}, and is stronger than pointwise convergence.

Topology

Definition 52.1 Topological Space^[?] (X, τ) : Is an ordered pair (X, τ) , where X is a set and τ is a topology^[def. 61.1] on X .

Definition 52.2 Topological Space^[?] (X, τ) : Is an ordered pair (X, τ) , where X is a set and τ is a topology^[def. 61.1] on X .

1. Weak Topologies

Definition 52.3 Weak Topology $\mathcal{C}(K; \mathbb{R})$: Is the coarsest topology s.t. all cont. linear functionals w.r.t. to the strong topology are continuous.
Neighbourhood Basis:

$$\{f || l_1 | < \epsilon_1, \dots, |l_n| < \epsilon_n, \forall \epsilon_i, \forall n, \forall \text{lin. functions } f\} \quad (52.1)$$

Note

The weak closure:

- is usually larger as the uniform closure, as for the weak closure there are many more convergence sequences
- is easier to calculate than the uniform closure

2. Compact Space

Corollary 52.1 Euclidean Space: In the euclidean case, a set $X \in \mathbb{R}$ is compact iff:

- it is closed^[def. 50.9]
- bounded

3. Closure

Definition 52.4 Closure of a Set^[?] $\text{cl}_{X, \tau}(S) / \bar{S}$: The closure of a subset S of a toplogical space^[def. 61.2] (X, τ) is defined equivalently by:

- is the union of S and its boundary ∂S .
- is the set S together with its limit points.

Note

If the topological space X, τ is clear from context, then the closure of a set S is often written simply as \bar{S} .

Corollary 52.2 Uniform Closure

$\|\cdot\|_\infty$: The uniform closure of a set of functions A is the space of all functions that can be approximated by a sequence (f_n) of uniformly-converging functions from A .^[def. 51.7] functions

Corollary 52.3 Weak Closure:

Logic

1. Boolean Algebra

1. Basic Operations

Definition 53.1 Conjunction/AND

\wedge :

Definition 53.2 Disjunction/OR

\vee :

Definition 53.3 Negation/NOT

\neg :

1.1.1. Expression as Integer

If the truth values {0, 1} are interpreted as integers then the basic operations can be represented with basic arithmetic operations.

$$\begin{aligned}x \wedge y &= xy = \min(x, y) \\x \vee y &= x + y = \max(x, y) \\-\bar{x} &= 1 - x \\x \oplus y &= (x + y) \cdot (\neg x + \neg y) = x \cdot \neg y + \neg x \cdot y\end{aligned}$$

Note: non-linearity of XOR

$$(x + y) \cdot (\neg x + \neg y) = -x^2 - y^2 - 2xy + 2x + 2y$$

2. Boolean Identities

Property 53.1 Idempotence:

$$x \wedge x \equiv x \quad \text{and} \quad x \vee x \equiv x \quad (53.1)$$

Property 53.2 Identity Laws:

$$x \wedge \text{true} \equiv x \quad \text{and} \quad x \vee \text{false} \equiv x \quad (53.2)$$

Property 53.3 Zero Law's:

$$x \wedge \text{false} \equiv \text{false} \quad \text{and} \quad x \vee \text{true} \equiv \text{true} \quad (53.3)$$

Property 53.4 Double Negation:

$$\neg\neg x \equiv x \quad (53.4)$$

Property 53.5 Complementation:

$$x \wedge \neg x \equiv \text{false} \quad \text{and} \quad x \vee \neg x \equiv \text{true} \quad (53.5)$$

Property 53.6 Commutativity:

$$x \vee y \equiv y \vee x \quad \text{and} \quad x \wedge y \equiv y \wedge x \quad (53.6)$$

Property 53.7 Associativity:

$$(x \vee y) \vee z \equiv x \vee (y \vee z) \quad (53.7)$$

$$(x \wedge y) \vee z \equiv x \wedge (y \wedge z) \quad (53.8)$$

Property 53.8 Distributivity:

$$x \vee (y \wedge z) \equiv (x \vee y) \wedge (x \vee z) \quad (53.9)$$

$$x \wedge (y \vee z) \equiv (x \wedge y) \vee (x \wedge z) \quad (53.10)$$

Property 53.9 De Morgan's Laws:

$$\neg(x \vee z) \equiv (\neg x \wedge \neg z) \quad (53.11)$$

$$\neg(x \wedge z) \equiv (\neg x \vee \neg z) \quad (53.12)$$

Note

The algebra axioms come in pairs that can be obtained by interchanging \wedge and \vee .

3. Normal Forms

Definition 53.4 Literal [example 53.1]:

Literals are atomic formulas or their negations

Definition 53.5 Negation Normal Form (NNF): A formula F is in negation normal form if the negation operator \neg is only applied to literals^[def. 53.4] and the only other operators are \wedge and \vee .

Definition 53.6 Conjunctive Normal Form (CNF): An boolean algebraic expression F is in CNF if it is a conjunction of clauses, where each clause is a disjunction of literals^[def. 53.4]

$L_{i,j}$:

$$F_{\text{CNF}} = \bigwedge_{i=1}^n \left(\bigvee_{j=1}^{m_i} L_{i,j} \right) \quad (53.13)$$

Definition 53.7 Disjunctive Normal Form (DNF): An boolean algebraic expression F is in DNF if it is a disjunction of clauses, where each clause is a conjunction of literals^[def. 53.4] $L_{i,j}$:

$$F_{\text{DNF}} = \bigvee_{i=1}^n \left(\bigwedge_{j=1}^{m_i} L_{i,j} \right) \quad (53.14)$$

Note

- true is a CNF with no clause and a single literal.
- false is a CNF with a single clause and no literals

1.3.1. Transformation to CNF and DNF

DNF

Algorithm 53.1:

- ① Using De Morgan's lawsProperty 53.9 and double negationProperty 53.4 transform F into Negation Normal Form^[def. 53.5]:

$$\begin{array}{ll} \neg\neg x & \text{by } x \\ \neg(x \wedge y) & \text{by } (\neg x \vee \neg y) \\ \neg(x \vee y) & \text{by } (\neg x \wedge \neg y) \\ \text{true} & \text{by } \text{false} \\ \text{false} & \text{by } \text{true} \end{array}$$

- ② Using distributive lawsProperty 53.8 substitute all:

$$\begin{array}{ll} x \wedge (y \vee z) & \text{by } (x \wedge y) \vee (x \wedge z) \\ (y \vee z) \wedge x & \text{by } (y \wedge x) \vee (z \wedge x) \\ x \wedge \text{true} & \text{by } \text{true} \\ \text{true} \wedge x & \text{by } \text{true} \end{array}$$

- ③ Using the identityProperty 53.2 and zero laws Property 53.3 remove true from any clause and delete all clauses containing false.

Note

For the CNF form simply use duality for step 2 and 3 i.e. swap \wedge and \vee and true and false.

Using Truth Tables [example 53.2]

To obtain a DNF formula from a truth table we need to have a conjunctive^[def. 53.3] for each row where F is true.

2. Examples

Example 53.1 Literals:

Boolean literals: $x, \neg y, s$

Not boolean literals: $\neg\neg x, (x \wedge y)$

Example 53.2 DNF from truth tables:

x	y	z	F
0	0	0	1
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	1
1	0	1	0
1	1	0	0
1	1	1	1

Need a conjunction of:

- $(\neg x \wedge \neg y \wedge \neg z)$
- $(\neg x \wedge y \wedge z)$
- $(x \wedge \neg y \wedge \neg z)$
- $(x \wedge y \wedge z)$

$(\neg x \wedge \neg y \wedge \neg z) \wedge (\neg x \wedge y \wedge z) \wedge (x \wedge \neg y \wedge \neg z) \wedge (x \wedge y \wedge z)$

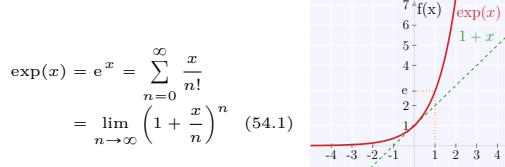
Calculus and Analysis

1. Functional Analysis

1. Elementary Functions

1.1.1. Exponential Numbers

Definition 54.1 Exponential Function



Definition 54.2 Exponential/Euler Number

$$e := \sum_{n=0}^{\infty} \frac{1}{n!} = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = 2.7182 \quad (54.2)$$

Properties Defining the Exponential Function

Property 54.1:

$$\exp(x+y) = \exp(x) + \exp(y) \quad (54.3)$$

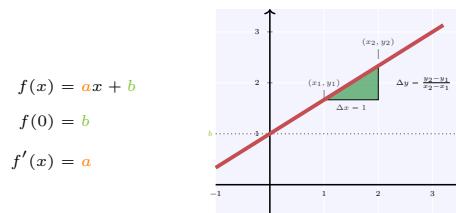
Property 54.2:

$$\exp(x) \leq 1+x \quad (54.4)$$

1.1.2. Affine Linear Functions

Definition 54.3 Affine Linear Function

$f(x) = ax + b$: An affine linear function are functions that can be defined by a scaling $s_a(x) = ax$ plus a translation $t_b(x) = x + b$:
 $M = \{f : \mathbb{R} \mapsto \mathbb{R} | f(x) = (s_a \circ t_b)(x) = ax + b, a, b \in \mathbb{R}\}$



Formula 54.1 [proof 54.1]

Linear Function from Point and slope

$f(x_0) = y_1$: Given a point (x_1, y_1) and a slope a we can derive:
 $f(x) = a \cdot (x - x_0) + y_0 = ax + (y_1 - ax_0)$

Formula 54.2 Linear Function from two Points:

$$f(x) = a \cdot (x - x_p) + y_p = ax + (y_p - ax_p) \quad (54.7)$$

$$a = \frac{y_1 - y_0}{x_1 - x_0} \quad p = \{0 \text{ or } 1\}$$

1.1.3. Polynomials

Definition 54.4 Polynomial:

A function $P_n : \mathbb{R} \mapsto \mathbb{R}$ is called *Polynomial*, if it can be represented in the form:

$$P_n(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1} + a_n x^n \quad (54.8)$$

Corollary 54.1 Degree n-of a Polynomial $\deg(P_n)$:

the degree of the polynomial is the highest exponent of the variable x , among all non-zero coefficients $a_i \neq 0$.

Definition 54.5 Monomial:

Is a polynomial with only one term.

Cubic Polynomials

Definition 54.6 Cubic Polynomials: Are polynomials of degree $\deg(P_n)$ 3 and have four coefficients:

$$f(x) = a_3 x^3 + a_2 x^2 + a_1 x + a_0 \quad (54.9)$$

2. Functional Compositions

Definition 54.7 Functional Compositions
 $f \circ g$: Let $f : A \mapsto B$ and $g : D \mapsto C$ be to mappings s.t. $\text{codom}(f) \subseteq D$ then we can define a composition function $(f \circ g) : A \mapsto C$ as:
 $h(x) = (g \circ f)(x) = g(f(x)) \quad \text{with } x \in A \quad (54.10)$

Corollary 54.2 Nested Functional Composition:

$$F_{k:1}(x) = (F_k \circ \dots \circ F_1)(x) = F_k(F_{k-1} \circ \dots \circ (F_1(x))) \quad (54.11)$$

2. Proofs

Proof 54.1 formula 54.1:

$$f(x_0) = y_0 = a x_0 + b \Rightarrow b = y_0 - a x_0$$

Theorem 54.1

First Fundamental Theorem of Calculus:

Let f be a continuous real-valued function defined on a closed interval $[a, b]$. Let F be the function defined $\forall x \in [a, b]$ by:

$$F(x) = \int_a^x f(t) dt \quad (54.12)$$

Then it follows:

$$F'(x) = f(x) \quad \forall x \in (a, b) \quad (54.13)$$

Theorem 54.2

Second Fundamental Theorem of Calculus:

Let f be a real-valued function on a closed interval $[a, b]$ and F an antiderivative of f in $[a, b]$: $F'(x) = f(x)$, then it follows if f is Riemann integrable on $[a, b]$:

$$\int_a^b f(t) dt = F(b) - F(a) \iff \int_a^x \frac{\partial}{\partial x} F(t) dt = F(x) \quad (54.14)$$

Definition 54.8 Domain of a function

Given a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, the set of all possible input values \mathcal{X} is called the domain of $f = \text{dom}(f)$.

Definition 54.9 Codomain/target set of a function

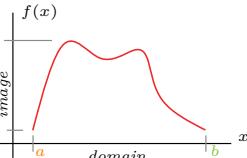
Given a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, the codomain of that function is the set \mathcal{Y} into which all of the output of the function is constrained to fall.

Definition 54.10 Image (Range) of a function: $f[\cdot]$

Given a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, the image of that function is the set to which the function can actually map:

$$\{y \in \mathcal{Y} | y = f(x), \forall x \in \mathcal{X}\} := f[\mathcal{X}] \quad (54.15)$$

Evaluating the function f at each element of a given subset A of its domain $\text{dom}(f)$ produces a set called the *image of A under (or through) f* . The image is thus a subset of a function's codomain.



Misnomer Range: The term Range is ambiguous s.t. certain books refer to it as codomain and other as image.

Definition 54.11 Inverse Image/Preimage $f^{-1}(\cdot)$:

Let $f : X \mapsto Y$ be a function, and A a subset set of its codomain Y .

Then the preimage of A under f is the set of all elements of the domain X , that map to elements in A under f :

$$f^{-1}(A) = \{x \in X : f(x) \subseteq A\} \quad (54.16)$$

Example 54.1 :

Given $f : \mathbb{R} \rightarrow \mathbb{R}$
defined by $f : x \mapsto x^2 \iff f(x) = x^2$
 $\text{dom}(f) = \mathbb{R}$, $\text{codom}(f) = \mathbb{R}$ but its image is $f[\mathbb{R}] = \mathbb{R}_+$.

Image (Range) of a subset

The image of a subset $A \subseteq \mathcal{X}$ under f is the subset $f[A] \subseteq \mathcal{Y}$ defined by:

$$f[A] = \{y \in \mathcal{Y} | y = f(x), \forall x \in A\} \quad (54.17)$$

Note: Range

The term range is ambiguous as it may refer to the image or the codomain, depending on the definition.
However, modern usage almost always uses range to mean image.

Definition 54.12 (strictly) Increasing Functions:

A function f is called monotonically increasing/increasing/non-decreasing if:
 $x \leq y \iff f(x) \leq f(y) \quad \forall x, y \in \text{dom}(f)$

$$(54.18)$$

And **strictly increasing** if:

$$x < y \iff f(x) < f(y) \quad \forall x, y \in \text{dom}(f) \quad (54.19)$$

Definition 54.13 (strictly) Decreasing Functions:

A function f is called monotonically decreasing/decreasing or non-increasing if:
 $x \geq y \iff f(x) \geq f(y) \quad \forall x, y \in \text{dom}(f)$

$$(54.20)$$

And **strictly decreasing** if:

$$x > y \iff f(x) > f(y) \quad \forall x, y \in \text{dom}(f) \quad (54.21)$$

Definition 54.14 Monotonic Function:

A function f is called monotonic iff either f is **increasing** or **decreasing**.

Definition 54.15 Linear Function:

A function $L : \mathbb{R}^n \mapsto \mathbb{R}^m$ is linear if and only if:

$$\begin{aligned} L(x+y) &= L(x) + L(y) \\ L(ax) &= aL(x) \end{aligned} \quad \forall x, y \in \mathbb{R}^n, a \in \mathbb{R}$$

Corollary 54.3 Linearity of Differentiation: The derivative of **any** linear combination of functions equals the same linear combination of the derivatives of the functions:

$$\frac{d}{dx} (a f(x) + b g(x)) = a \frac{d}{dx} f(x) + b \frac{d}{dx} g(x) \quad a, b \in \mathbb{R} \quad (54.22)$$

Definition 54.16 Quadratic Function:

A function $f : \mathbb{R}^n \mapsto \mathbb{R}^m$ is quadratic if it can be written in the form:

$$f(x) = \frac{1}{2} x^T A x + b^T x + c \quad (54.23)$$

3. Norms

1. Infinity/Supremum Norm

Definition 54.17 Infinity/Supremum Norm:

$$\|f\|_\infty := \sup_{x \in \text{dom}(f)} |f(x)| \quad (54.24)$$

Note

In order to make this a proper norm one usually considers **bounded functions** s.t.:

$$\|f\|_\infty \leq M < \infty$$

Corollary 54.4 Ininity Norm induced Metric: The infinity norm naturally induces a metric $d := (f, g) := \|f - g\|_\infty$

$$(54.25)$$

4. Smoothness

Definition 54.18 Smoothness of a Function \mathcal{C}^k :

Given a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, the function is said to be of class k if it is differentiable up to order k and continuous, on its entire domain:

$$f \in \mathcal{C}^k(\mathcal{X}) \iff \exists f', f'', \dots, f^{(k)} \text{ continuous} \quad (54.26)$$

Note

- P.w. continuous \neq continuous.
- A function of that is k times differentiable must at least be of class \mathcal{C}^{k-1} .
- $\mathcal{C}^m(\mathcal{X}) \subset \mathcal{C}^{m-1}, \dots, \mathcal{C}^1 \subset \mathcal{C}^0$
- Continuity is implied by the differentiability of all derivatives of up to order $k-1$.

4.0.1. Continuous Functions

Definition 54.19 Continuous Function \mathcal{C}^0 : Functions that do not have any jumps or peaks.

4.0.2. Piece wise Continuous Functions

Definition 54.20 Piecewise Linear Functions

$$\mathcal{C}_{pw}^0$$

4.0.3. Continuously Differentiable Function

Corollary 54.5 Continuously Differentiable Function: C^1 : Is the class of functions that consists of all differentiable functions whose derivative is continuous.

Hence a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ of the class must satisfy:

$$f \in C^1(\mathcal{X}) \iff f' \text{ continuous} \quad (54.27)$$

4.0.4. Smooth Functions

Corollary 54.6 Smooth Function C^∞ : Is a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ that has derivatives infinitely many times differentiable.

$$f \in C^\infty(\mathcal{X}) \iff f', f'', \dots, f^{(\infty)} \quad (54.28)$$

1. Lipschitz Continuous Functions

Often functions are not differentiable but we still want to state something about the rate of change of a function \Rightarrow hence we need a weaker notion of differentiability.

Definition 54.21 Lipschitz Continuity:

A Lipschitz continuous function is a function f whose rate of change is bound by a Lipschitz Constant L :

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \quad L > 0 \quad (54.29)$$

Note

This property is useful as it allows us to conclude that a small perturbation in the input (i.e. of an algorithm) will result in small changes of the output \Rightarrow tells us something about robustness.

4.1.1. Lipschitz Continuous Gradient

Definition 54.22 Lipschitz Continuous Gradient:

A continuously differentiable function $f : \mathbb{R}^d \mapsto \mathbb{R}$ has L -Lipschitz continuous gradient if it satisfies:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), \quad L > 0 \quad (54.30)$$

if $f \in C^2$, this is equivalent to:

$$\nabla^2 f(\mathbf{x}) \leq L \mathbf{I} \quad \forall \mathbf{x} \in \text{dom}(f), \quad L > 0 \quad (54.31)$$

Lemma 54.1 Descent Lemma

[Poorfs 54.5,??]:

If a function $f : \mathbb{R}^d \mapsto \mathbb{R}$ has Lipschitz continuous gradient eq. (54.30) over its domain, then it holds that:

$$|f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (54.32)$$

Note

If f is twice differentiable then the largest eigenvalue of the Hessian (Definition 55.8) of f is uniformly upper bounded by L

2. L-Smooth Functions

Definition 54.23 L-Smoothness:

A L -smooth function is a function $f : \mathbb{R}^d \mapsto \mathbb{R}$ that satisfies:

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

with $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), \quad L > 0 \quad (54.33)$

If f is a twice differentiable this is equivalent to:

$$\nabla^2 f(\mathbf{x}) \leq L \mathbf{I} \quad L > 0 \quad (54.34)$$

Theorem 54.3 [proof 54.6]

L-Smoothness of convex functions:

A convex and L-Smooth function (^{def. 54.23}) has a Lipschitz continuous gradient eq. (54.30) thus it holds that:

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (54.35)$$

Note

L -smoothness is a weaker condition than L -Lipschitz continuous gradients

5. Convexity and Concavity

Definition 54.24 Convex Functions:

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if it satisfies:

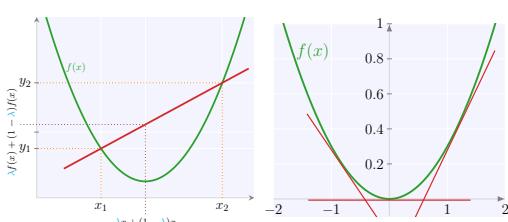
$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \lambda \in [0, 1] \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (54.36)$$

If f is a differentiable function this is equivalent to:

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (54.37)$$

If f is a twice differentiable function this is equivalent to:

$$\nabla^2 f(\mathbf{x}) \geq 0 \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (54.38)$$



Definition 54.25 Concave Functions:

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is concave if it satisfies:

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \geq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \lambda \in [0, 1] \quad (54.39)$$

Corollary 54.7 Convexity \rightarrow global minimina: Convexity implies that all local minima (if they exist) are global minima.

1. Properties

Property 54.3 Monotonicity of the Derivative:

convex	$f'(\mathbf{a}) < f'(\mathbf{b})$
If $f : \mathbb{R} \mapsto \mathbb{R}$ is	$\mathbf{a} < \mathbf{b}, \quad \mathbf{a}, \mathbf{b} \in \mathbb{R}$
concave	$f'(\mathbf{a}) > f'(\mathbf{b})$

$$(54.40)$$

1.1. Properties that preserve convexity

Property 54.4 Non-negative weighted Sums: Let f be a convex function then $g(\mathbf{x})$ is convex as well:

$$g(\mathbf{x}) = \sum_{i=1}^n \alpha_i f_i(\mathbf{x}) \quad \forall \alpha_i > 0$$

Property 54.5 Composition of Affine Mappings: Let f be a convex function then $g(\mathbf{x})$ is convex as well:

$$g(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$$

Property 54.6 Pointwise Maxima: Let f be a convex function then $g(\mathbf{x})$ is convex as well:

$$g(\mathbf{x}) = \max_i \{f_i(\mathbf{x})\}$$

2. Strict Convexity/Concavity

Definition 54.26 Strictly Convex Functions:

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is strictly convex if it satisfies:

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) < \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \lambda \in [0, 1] \quad (54.41)$$

If f is a differentiable function this is equivalent to:

$$f(\mathbf{x}) > f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (54.41)$$

If f is a twice differentiable function this is equivalent to:

$$\nabla^2 f(\mathbf{x}) > 0 \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (54.42)$$

Intuition

- Convexity implies that a function f is bound by/below a linear interpolation from x to y and strong convexity that f is strictly bound/below.
- eq. (54.41) implies that $f(\mathbf{x})$ is above the tangent $f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
- ?? implies that $f(\mathbf{x})$ is flat or curved upwards

Corollary 54.8 Strict Convexity \rightarrow Uniqueness:

Strict convexity implies a unique minimizer \iff at most one global minimum.

Corollary 54.9 : A twice differentiable function of one variable $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex on an interval $\mathcal{X} = [\mathbf{a}, \mathbf{b}]$ if and only if its second derivative is non-negative on that interval \mathcal{X} :

$$f''(\mathbf{x}) \geq 0 \quad \forall \mathbf{x} \in \mathcal{X} \quad (54.43)$$

3. Strong Convexity/Concavity

Definition 54.27 μ -Strong Convexity:

Let \mathcal{X} be a Banach space over $K = \mathbb{R}, \mathbb{C}$. A function $f : \mathcal{X} \rightarrow \mathbb{R}$ is called strongly convex iff the following equation holds:

$$f(t\mathbf{x} + (1 - t)\mathbf{y}) \leq t f(\mathbf{x}) + (1 - t)f(\mathbf{y}) - \frac{\mu(1 - t)}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}, \quad t \in [0, 1], \quad \mu > 0$$

If $f \in C^1 \iff f$ is differentiable, this is equivalent to:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (54.44)$$

If $f \in C^2 \iff f$ is twice differentiable, this is equivalent to:

$$\nabla^2 f(\mathbf{x}) \geq \mu \mathbf{I} \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X} \quad \mu > 0 \quad (54.45)$$

Corollary 54.10

Strong Convexity implies Strict Convexity:

Property 54.7:

$$f(\mathbf{y}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) + \frac{1}{2\mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \quad (54.46)$$

Intuition

Strong convexity implies that a function f is lower bounded by its second order (quadratic) approximation, rather than only its first order (linear) approximation.

Size of μ

The parameter μ specifies how strongly the bounding quadratic function/approximation is.

Proof 54.2: eq. (54.45) analogously to Proof eq. (54.34)

Note

If f is twice differentiable then the smallest eigenvalue of the Hessian (^{def. 55.8}) of f is uniformly lower bounded by μ . Hence strong convexity can be considered as the analogous to smoothness

Example 54.2 Quadratic Function: A quadratic function eq. (54.23) is convex if:

$$\nabla_{\mathbf{x}}^2 \text{eq. (54.23)} = \mathbf{A} \geq 0 \quad (54.47)$$

Corollary 54.11 :

Strong convexity \Rightarrow Strict convexity \Rightarrow Convexity

Functions

Even Functions:

have rotational symmetry with respect to the origin.

\Rightarrow Geometrically: its graph remains unchanged after reflection about the y-axis.

$$f(-x) = f(x) \quad (54.48)$$

Odd Functions: are symmetric w.r.t. to the y-axis.

\Rightarrow Geometrically: its graph remains unchanged after rotation of 180 degrees about the origin.

$$f(-x) = -f(x) \quad (54.49)$$

Examples

Even: $\cos x, |x|, c, x^2, x^4, \dots, \exp(-x^2/2)$.

Odd: $\sin x, \tan x, x, x^3, x^5, \dots$

x-Shift: $f(x - c) \Rightarrow$ shift to the right

$$f(x + c) \Rightarrow$$
 shift to the left (54.50)

y-Shift: $f(x) \pm c \Rightarrow$ shift up/down (54.51)

Proof 54.3: eq. (54.50) $f(x_n - c)$ we take the x -value at x_n but take the y -value at $x_0 : x_n = x_0 - c$

\Rightarrow we shift the function to x_n .

Euler's formula

$$e^{\pm ix} = \cos x \pm i \sin x \quad (54.52)$$

Euler's Identity

$$e^{\pm i\pi} = -1 \quad (54.53)$$

Note

$$e^n = 1 \Leftrightarrow n = i 2\pi k, \quad k \in \mathbb{N} \quad (54.54)$$

Corollary 54.12 Every norm is a convex function: By using definition ^{def. 54.24} and the triangular inequality it follows (with the exception of the L0-norm):

$$\|\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}\| \leq \lambda \|\mathbf{x}\| + (1 - \lambda) \|\mathbf{y}\|$$

4. Taylor Expansion

Definition 54.28 Taylor Expansion:

$$T_n(\mathbf{x}) = \sum_{i=0}^n \frac{1}{i!} f^{(i)}(\mathbf{x}_0) \cdot (\mathbf{x} - \mathbf{x}_0)^{(i)} \quad (54.55)$$

$$= f(\mathbf{x}_0) + f'(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} f''(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)^2 + \mathcal{O}(x^3) \quad (54.56)$$

Note

If we chose Δx small enough it is sufficient to look only at the first two terms.

Definition 54.29 Incremental Taylor:

Goal: evaluate $T_n(\mathbf{x})$ (eq. (54.56)) at the point $\mathbf{x}_0 + \Delta x$ in order to propagate the function $f(\mathbf{x})$ by $h = \Delta x$:

$$T_n(\mathbf{x}_0 \pm h) = \sum_{i=0}^n \frac{h^i}{i!} f^{(i)}(\mathbf{x}_0) \mathbf{i}^{-1} \quad (54.57)$$

$$= f(\mathbf{x}_0) \pm h f'(\mathbf{x}_0) + \frac{h^2}{2} f''(\mathbf{x}_0) \pm f'''(\mathbf{x}_0) h^3 + \mathcal{O}(h^4) \quad (54.58)$$

Note

If we chose Δx small enough it is sufficient to look only at the first two terms.

Definition 54.30 Multidimensional Taylor:

Suppose $X \in \mathbb{R}^n$ is open, $\mathbf{x} \in X$, $f : X \mapsto \mathbb{R}$ and $f \in C^2$ then it holds that

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + \nabla_{\mathbf{x}} f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top H(\mathbf{x} - \mathbf{x}_0) \quad (54.58)$$

Definition 54.31 Argmax:

The argmax of a function defined on a set D is given by:

$$\arg \max_{\mathbf{x} \in D} f(\mathbf{x}) = \{ \mathbf{x} | f(\mathbf{x}) \geq f(\mathbf{y}), \forall \mathbf{y} \in D \} \quad (54.59)$$

Definition 54.32 Argmin: The argmin of a function defined on a set D is given by:

$$\arg \min_{x \in D} f(x) = \{x | f(x) \leq f(y), \forall y \in D\} \quad (54.60)$$

Note

The supremum/infinum is necessary to handle unbound function that seem to converge and for which the max/min does not exist as the argmax/argmin may be empty.

E.g. consider $-e^x/e^x$ for which the max/min converges toward 0 but will never reach s.t. we can always choose a bigger $x \Rightarrow$ there exists no argmax/argmin \Rightarrow need to bound the functions from above/below \iff infimum/supremum.

Corollary 54.13 Relationship $\arg \min \leftrightarrow \arg \max$:

$$\arg \min_{x \in D} f(x) = \arg \max_{x \in D} -f(x) \quad (54.61)$$

Property 54.8 Argmax Identities:

1. **Shifting:** $\forall \lambda \text{ const } \arg \max f(x) = \arg \max f(x) + \lambda \quad (54.62)$

2. **Positive Scaling:** $\forall \lambda > 0 \text{ const } \arg \max f(x) = \arg \max \lambda f(x) \quad (54.63)$

3. **Negative Scaling:** $\forall \lambda < 0 \text{ const } \arg \max f(x) = \arg \min \lambda f(x) \quad (54.64)$

4. **Positive Functions:** $\forall \arg \max f(x) > 0, \forall x \in \text{dom}(f)$

$$\arg \max f(x) = \arg \min \frac{1}{f(x)} \quad (54.65)$$

5. **Strictly Monotonic Functions:** for all strictly monotonic increasing functions^[def. 54.12] g it holds that:

$$\arg \max g(f(x)) = \arg \max f(x) \quad (54.66)$$

Definition 54.33 Max: The maximum of a function f defined on the set D is given by:

$$\max_{x \in D} f(x) = f(x^*) \quad \text{with} \quad \forall x^* \in \arg \max_{x \in D} f(x) \quad (54.67)$$

Definition 54.34 Min: The minimum of a function f defined on the set D is given by:

$$\min_{x \in D} f(x) = f(x^*) \quad \text{with} \quad \forall x^* \in \arg \min_{x \in D} f(x) \quad (54.68)$$

Corollary 54.14 Relationship $\min \leftrightarrow \max$:

$$\min_{x \in D} f(x) = -\max_{x \in D} -f(x) \quad (54.69)$$

Property 54.9 Max Identities:

1. **Shifting:** $\forall \lambda \text{ const } \max \{f(x) + \lambda\} = \lambda + \max f(x) \quad (54.70)$

2. **Positive Scaling:** $\forall \lambda > 0 \text{ const } \max \lambda f(x) = \lambda \max f(x) \quad (54.71)$

3. **Negative Scaling:** $\forall \lambda < 0 \text{ const } \max \lambda f(x) = \lambda \min f(x) \quad (54.72)$

4. **Positive Functions:** $\forall \arg \max f(x) > 0, \forall x \in \text{dom}(f)$

$$\max_{x \in D} \frac{1}{f(x)} = \frac{1}{\min f(x)} \quad (54.73)$$

5. **Strictly Monotonic Functions:** for all strictly monotonic increasing functions^[def. 54.12] g it holds that:

$$\max g(f(x)) = g(\max f(x)) \quad (54.74)$$

Definition 54.35 Supremum: The supremum of a function defined on a set D is given by:

$$\sup_{x \in D} f(x) = \{y | y \geq f(x), \forall x \in D\} = \min_{y | y \geq f(x), \forall x \in D} y \quad (54.75)$$

and is the smallest value y that is equal or greater $f(x)$ for any $x \iff$ smallest upper bound.

Definition 54.36 Infimum: The infimum of a function defined on a set D is given by:

$$\inf_{x \in D} f(x) = \{y | y \leq f(x), \forall x \in D\} = \max_{y | y \leq f(x), \forall x \in D} y \quad (54.76)$$

and is the biggest value y that is equal or smaller $f(x)$ for any $x \iff$ largest lower bound.

Corollary 54.15 Relationship $\sup \leftrightarrow \inf$:

$$\sup_{x \in D} f(x) = -\inf_{x \in D} -f(x) \quad (54.77)$$

Proof 54.5: ?? for C^2 functions:

$$f(\mathbf{y}) \stackrel{\text{Taylor}}{=} f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top \nabla^2 f(\mathbf{z})(\mathbf{y} - \mathbf{x})$$

Now we plug in $\nabla^2 f(\mathbf{x})$ and recover eq. (54.33):

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top L(\mathbf{y} - \mathbf{x})$$

Proof 54.6: theorem 54.3:

With the definition of convexity for a differentiable function (eq. (54.41)) it follows

$$f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \geq 0$$

$$\Rightarrow |f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})|$$

$$\text{if eq. (54.41)} \quad f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})$$

with lemma 54.1 and [def. 54.23] it follows theorem 54.3

Definition 54.37 Time-invariant system (TIS): A function f is called time-invariant, if shifting the input in time leads to the same output shifted in time by the same amount.

$$y(t) = f(x(t), t) \xrightarrow[\forall \tau]{\text{time-invariance}} y(t - \tau) = f(x(t - \tau), t) \quad (54.78)$$

Definition 54.38 Inverse Function $g = f^{-1}$:

A function g is the inverse function of the function $f : A \subset \mathbb{R} \rightarrow B \subset \mathbb{R}$ if

$$f(g(x)) = x \quad \forall x \in \text{dom}(g) \quad (54.79)$$

and

$$g(f(u)) = u \quad \forall u \in \text{dom}(f) \quad (54.80)$$

Property 54.10

Reflective Property of Inverse Functions: f contains (a, b) if and only if f^{-1} contains (b, a) .
The line $y = x$ is a symmetry line for f and f^{-1} .

Theorem 54.5 The Existence of an Inverse Function:

A function has an inverse function if and only if it is one-to-one.

Corollary 54.16 Inverse functions and strict monotonicity: If a function f is strictly monotonic^[def. 54.14] on its entire domain, then it is one-to-one and therefore has an inverse function.

6. Special Functions

1. The Gamma Function

Definition 54.39 The gamma function $\Gamma(\alpha)$: Is extension of the factorial function ?? to the real and complex numbers (with a positive real part):

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx \quad \Re(z) > 0 \quad (54.81)$$

$$\Gamma(n) \stackrel{n \in \mathbb{N}}{\iff} \Gamma(n) = (n-1)!$$

7. Proofs

Proof 54.4: lemma 54.1 for C^1 functions:

Let $g(t) \equiv f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))$ from the FToC (theorem 54.2) we know that:

$$\int_0^1 g'(t) dt = g(1) - g(0) = f(\mathbf{x}) - f(\mathbf{y})$$

It then follows from the reverse:

$$\begin{aligned} & |f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})| \\ \stackrel{\text{Chain. R}}{\stackrel{\text{FToC}}{=}} & \left| \int_0^1 \nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^\top (\mathbf{x} - \mathbf{y}) dt - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \right| \\ = & \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^\top (\mathbf{x} - \mathbf{y}) dt \right| \\ = & \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^\top (\mathbf{x} - \mathbf{y}) dt \right| \\ \stackrel{\text{C.S.}}{\leq} & \int_0^1 \| \nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}) \| \cdot \| \mathbf{x} - \mathbf{y} \| dt \\ \stackrel{\text{eq. (54.30)}}{=} & \int_0^1 L \| \mathbf{y} + t(\mathbf{x} - \mathbf{y}) - \mathbf{y} \| \cdot \| \mathbf{x} - \mathbf{y} \| dt \\ = & L \| \mathbf{x} - \mathbf{y} \|^2 \int_0^1 t dt \\ = & L \| \mathbf{x} - \mathbf{y} \|^2 \left[\frac{t^2}{2} \right]_0^1 = \frac{L}{2} \| \mathbf{x} - \mathbf{y} \|^2 \end{aligned}$$

Differential Calculus

1. Mean Value Theorem

Theorem 55.1 Mean Value Theorem: Let $f : [a, b] \rightarrow \mathbb{R}$ be continuous function, differentiable on the open interval (a, b) , with $a < b$. Then there exist some $c \in (a, b)$ s.t.

$$f'(c) = \frac{f(b) - f(a)}{b - a} = \frac{1}{b - a} \int_a^b f(x) dx \quad (55.1)$$

2. The Product Rule

Rule 55.1 (Product /Leibniz Rule).

Let u, v be two differentiable functions $u, v \in C^1$ then it holds that:

$$\frac{d(u(x)v(x))}{dx} = (uv)' = u'v + v'u \quad (55.2)$$

3. The Chain Rule

Formula 55.1 Generalized Chain Rule:

Let $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and $\mathbf{G} : \mathbb{R}^k \rightarrow \mathbb{R}^m$ be to general maps then it holds:

$$\begin{aligned} \frac{\partial(\mathbf{G} \circ \mathbf{F})}{\partial \mathbf{x}} &= \underbrace{(\partial \mathbf{G} \circ \mathbf{F})}_{\mathbb{R}^n \mapsto \mathbb{R}^{m \times n}} \cdot \underbrace{\partial \mathbf{F}}_{\mathbb{R}^n \mapsto \mathbb{R}^{k \times n}} \\ &\quad \partial \mathbf{G} : \mathbb{R}^k \mapsto \mathbb{R}^{m \times k} \end{aligned} \quad (55.3)$$

4. Directional Derivative

5. Partial Differentiation

Definition 55.1 Partial Derivative:

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a real valued function, its partial derivative $\partial_i f : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined as the directional derivative?? along the coordinate axis of one of its variables:

$$\begin{aligned} \partial_i f(\mathbf{x}) &= D_{x_i} f = \lim_{h \rightarrow 0} \frac{f(\mathbf{x}, x_i \leftarrow x_i + h) - f(\mathbf{x})}{h} \\ &= \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{h} \end{aligned} \quad (55.4)$$

1. The Gradient

5.1.1. The Nabla Operator

Definition 55.2 Nabla Operator/Del

∇ : Given a cartesian coordinate system \mathbb{R}^n with coordinates x_1, \dots, x_n and associated unit vectors $\hat{\mathbf{e}}_1, \dots, \hat{\mathbf{e}}_n$ its del operator is defined as:

$$\nabla = \sum_{i=1}^n \frac{\partial}{\partial x_i} \hat{\mathbf{e}}_i = \begin{bmatrix} \frac{\partial}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_n}(\mathbf{x}) \end{bmatrix} \quad (55.5)$$

Definition 55.3 Gradient:

Given a scalar valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ its gradient $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is defined as vector \mathbb{R}^n of the partial derivatives w.r.t. all coordinate axes:

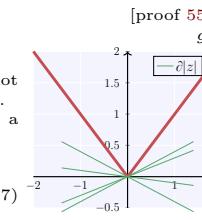
$$\text{grad } f(\mathbf{x}) := \nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \frac{\partial f}{\partial x_2}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix} = \left(\frac{\partial f}{\partial \mathbf{x}} \right)^\top \quad (55.6)$$

5.1.2. The Subderivative

Definition 55.4 Subgradient

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous (not necessarily differentiable) function. $\mathbf{g} \in \mathbb{R}^n$ is a subgradient of f at a point $\mathbf{x}_0 \in \mathbb{R}^n$ if it satisfies:

$$g : f(\mathbf{x}) - f(\mathbf{x}_0) \geq g^\top(\mathbf{x} - \mathbf{x}_0) \quad (55.7)$$



1. The Hessian

Definition 55.8 Hessian Matrix:

Given a function $f : \mathbb{R} \mapsto \mathbb{R}^n$ its Hessian $\mathbb{R}^{n \times n}$ is defined as:

$$\mathbf{H}(\mathbf{f})(\mathbf{x}) = \mathbf{H}_f(\mathbf{x}) = \mathbf{J}(\nabla \mathbf{f}(\mathbf{x}))^\top \quad (55.10)$$

$$= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) & & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{x}) \\ \vdots & & \ddots & \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{bmatrix}$$

and it corresponds to the Jacobian of the Gradient.

Due to the differentiability and theorem 55.2 it follows that the Hessian is (if it exists):

- Symmetric
- Real

Corollary 55.2 Eigenvector basis of the Hessian: Due to the fact that the Hessian is real and symmetric we can decompose it into a set of real eigenvalues and an orthogonal basis of eigenvectors $\{(\lambda_1, \mathbf{v}_1), \dots, (\lambda_n, \mathbf{v}_n)\}$.

Not let \mathbf{d} be a directional unit vector then the second derivative in that direction is given by:

$$\mathbf{d}^\top \mathbf{H} \mathbf{d} \iff \mathbf{d}^\top \sum_{i=1}^n \lambda_i \mathbf{v}_i \iff \mathbf{d}^\top \lambda_j \mathbf{v}_j$$

- The eigenvectors that have smaller angle with \mathbf{d} have bigger weight/eigenvalues
- The minimum/maximum eigenvalue determines the minimum/maximum second derivative

2. The Jacobian

Definition 55.6 Jacobian/Jacobi Matrix

$\mathbf{Df}, \mathbf{J_f}$:

Given a vector valued function

$\mathbf{f} : \mathbb{R}^n \mapsto \mathbb{R}^m$ its derivative $\mathbf{J}_f : \mathbb{R}^n \mapsto \mathbb{R}^{m \times n}$

with components $\partial_{ij} \mathbf{f} = \partial_i f_j : \mathbb{R}^n \mapsto \mathbb{R}$ is a vector valued function defined as:

$$\mathbf{J}(\mathbf{f}(\mathbf{x})) = \mathbf{J}_f(\mathbf{x}) = \mathbf{Df} = \frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}) = \frac{\partial(f_1, \dots, f_m)}{\partial(x_1, \dots, x_n)}(\mathbf{x}) \quad (55.9)$$

$$\begin{aligned} &= \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix} \\ &= \begin{bmatrix} \nabla^\top f_1 \\ \vdots \\ \nabla^\top f_m \end{bmatrix} \end{aligned}$$

Explanation 55.1: Rows of the Jacobian are transposed gradients^[def. 55.3] of the component functions f_1, \dots, f_m .

Corollary 55.1 :

6. Second Order Derivatives

Definition 55.7 Second Order Derivative $\frac{\partial^2}{\partial x_i \partial x_j}$:

Theorem 55.2

Symmetry of second derivatives/Schwartz's Theorem:

Given a continuous and twice differentiable function $f : \mathbb{R}^n \mapsto \mathbb{R}$ then its second order partial derivatives commute:

$$\frac{\partial}{\partial x_i} \frac{\partial f}{\partial x_j} = \frac{\partial}{\partial x_j} \frac{\partial f}{\partial x_i}$$

Corollary 55.4 Second Derivative Test $f : \mathbb{R} \mapsto \mathbb{R}$:

Suppose $f : \mathbb{R} \mapsto \mathbb{R}$ is twice differentiable at a stationary point \mathbf{x} ^[def. 55.9] then it follows that:

- $f''(x + \epsilon) > 0 \Rightarrow$ slope points uphill
- $f''(x - \epsilon) < 0 \Rightarrow$ slope points downhill
- $f(x)$ is a local minimum
- $f''(x + \epsilon) < 0 \Rightarrow$ slope points downhill
- $f''(x - \epsilon) > 0 \Rightarrow$ slope points uphill
- $f(x)$ is a local maximum

$\epsilon > 0$ sufficiently small enough

Corollary 55.4 Second Derivative Test $f : \mathbb{R}^n \mapsto \mathbb{R}$:

Suppose $f : \mathbb{R}^n \mapsto \mathbb{R}$ is twice differentiable at a stationary point \mathbf{x} ^[def. 55.9] then it follows that:

- If \mathbf{H} is p.d $\iff \forall \lambda_i > 0 \in \mathbf{H} \rightarrow f(\mathbf{x})$ is a local min.
- If \mathbf{H} is n.d $\iff \forall \lambda_i < 0 \in \mathbf{H} \rightarrow f(\mathbf{x})$ is a local max.
- If $\exists \lambda_i > 0 \in \mathbf{H}$ and $\exists \lambda_j < 0 \in \mathbf{H}$ then \mathbf{x} is a local maximum in one cross section of f but a local minimum in another
- If $\exists \lambda_i = 0 \in \mathbf{H}$ and all other eigenvalues have the same sign the test is inconclusive as it is in the cross section corresponding to the zero eigenvalue.

Note

If \mathbf{H} is positive definite for a minima \mathbf{x}^* of a quadratic function f then this point must be a global minimum of that function.

8. Proofs

Proof 55.1: Definition 55.4 $f(\mathbf{x}) \geq f(\mathbf{x}_0) + \mathbf{g}^\top(\mathbf{x} - \mathbf{x}_0) \quad \forall \mathbf{x} \in \mathbb{R}^n$ corresponds to a line (see formula 54.1) at the point \mathbf{x}_0 with slope \mathbf{g}^\top .

Thus we search for all lines with smaller slope then function graph.

9. Examples

Example 55.1 Subderivatives Absolute Value Function
 $|x|: f : \mathbb{R} \mapsto \mathbb{R}$ with $f(x) = |x|$ at the point $x = 0$ it holds:
 $f(x) - f(0) \geq g x \implies$ the interval $[-1; 1]$

For $x \neq 0$ the subgradient is equal to the gradient. Thus it follows for the subderivatives/differentials:

$$\partial|x| = \begin{cases} -1 & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Integral Calculus

Theorem 56.1 Important Integral Properties:

Addition $\int_a^b f(x) dx = \int_a^c f(x) dx + \int_c^b f(x) dx \quad (56.1)$

Reflection $\int_a^b f(x) dx = - \int_b^a f(x) dx \quad (56.2)$

Translation $\int_a^b f(x) dx \stackrel{u:=x\pm c}{=} \int_{a\pm c}^{b\pm c} f(x \mp c) dx \quad (56.3)$

f Odd $\int_{-a}^a f(x) dx = 0 \quad (56.4)$

f Even $\int_{-a}^a f(x) dx = 2 \int_0^a f(x) dx \quad (56.5)$

Proof 56.1: eqs. (56.4) and (56.5)

$$\begin{aligned} I := \int_{-a}^a f(x) dx &= \int_{-a}^0 f(x) dx + \int_0^a f(x) dx \\ dt = -dx &\quad \int_0^0 f(-x) dx + \int_0^a f(x) dx \\ &= \int_0^a f(-x) + f(x) dx = \begin{cases} 0 & \text{if } f \text{ odd} \\ 2I & \text{if } f \text{ even} \end{cases} \end{aligned}$$

Definition 56.1 Integration by Parts:

$$\int_a^b u dv = uv \Big|_a^b - \int_a^b v du \quad (56.6)$$

1. Integral Theorems

1. Greens Identities

Theorem 56.2 Greens First Identity:

Let $\bar{\Omega} = \Omega \cup \partial\Omega$, for all vector fields $j \in (\mathcal{C}_{pw}^1(\bar{\Omega}))^d$ and scalar functions $v \in \mathcal{C}_{pw}^1(\bar{\Omega})$ it holds:

$$\int_{\Omega} j^T \operatorname{grad} v dx = - \int_{\Omega} \operatorname{div} j v dx + \int_{\partial\Omega} j^T n v dS \quad (56.7)$$

Differential Equations

Definition 56.2

[??]

Differential Operator:

A differential operator \mathcal{L} is a mapping of a suitable function space onto another function space, involving only values of the function argument and its derivatives in the same point:
 $\mathcal{L} : C^n(\Omega) \rightarrow C^k(\Omega)$, $k < n$

Note: \mathcal{L} is a differential operator of order $k - n$.

Definition 56.3 Linear Differential Operator:

Is a differential operator \mathcal{L} that satisfies:

$$\mathcal{L}(\alpha u + \beta v) = \alpha \mathcal{L}(u) + \beta \mathcal{L}(v) \quad \forall \alpha, \beta \in \mathbb{R} \quad (56.8)$$

Ordinary Differential Quations

Partial Differential Equations (PDE)s

Definition 58.1 Partial Differential Equation:

Let $\mathbf{u} = \mathbf{u}(x_1, \dots, x_n) : \mathbb{R}^k \rightarrow \mathbb{R}^l$ be an unknown function depending on $\mathbf{x} = (x_1, \dots, x_k)$ and let f be a known function. The known function F , depending on differentials of the unknown function \mathbf{u} is called a Partial Differential equation:

$$F\left(\mathbf{u}, \frac{\partial \mathbf{u}}{\partial x_1}, \dots, \frac{\partial \mathbf{u}}{\partial x_i}, \dots, \frac{\partial \mathbf{u}}{\partial x_j}, f\right) = F(\mathbf{u}, D\mathbf{u}, \dots, D^n \mathbf{u}, f) = 0$$

or
 $\mathcal{L}(\mathbf{u}) = f \quad \text{in } \Omega \quad (58.1)$

Corollary 58.1 Dependent Variables:

$$\mathbf{u} : \mathbb{R}^k \rightarrow \mathbb{R}^l \quad (58.2)$$

Corollary 58.2 Independent Variables:

$$\mathbf{x} = (x_1, \dots, x_k) \quad (58.3)$$

Definition 58.2 Order

n:

Is the highest partial derivative that appears in a PDE.

1. Algebraic Types

1. Linearity

Definition 58.3

[??]

Linear PDEs:

A linear PDE naturally defines a linear operator [def. 56.3]. A linear PDE must be linear regarding the unknown function \mathbf{u} . In other words all dependent variables \mathbf{u} and their corresponding derivatives depend only on the independent variables x_1, x_2, \dots, x_m :

$$a(x, y)\mathbf{u}_x + b(x, y)\mathbf{u}_y + c(x, y)\mathbf{u} = d(x, y) \quad (58.4)$$

Definition 58.4

[??]

Semilinear PDEs:

Are PDEs whose coefficients of the highest order n -terms are functions depending only on the independent variables but not onto the dependent variables \mathbf{u} or their derivatives.

Thus the PDE is linear regarding to the highest order terms:

$$a(x, y)\mathbf{u}_x + b(x, y)\mathbf{u}_y = c(x, y, \mathbf{u}) \quad (58.5)$$

Definition 58.5

[??]

Quasilinear PDEs:

Are PDEs whose coefficients of the highest order (n) terms are functions only depending on the independent variables and on the dependent variables \mathbf{u} and their derivatives up to an order $m < n$, that is smaller than the highest order terms n :

$$a(x, y, \mathbf{u})\mathbf{u}_x + b(x, y, \mathbf{u})\mathbf{u}_y = c(x, y, \mathbf{u}) \quad (58.6)$$

Definition 58.6

[??]

Fully Non-linear PDEs:

Are PDEs where all terms of the highest order n are non-linear:

$$a(x, y, \mathbf{u}, \mathbf{u}')\mathbf{u}_x + b(x, y, \mathbf{u}, \mathbf{u}')\mathbf{u}_y = c(x, y, \mathbf{u}) \quad (58.7)$$

Note: $\neg(\text{Quasilinear} \Leftrightarrow \text{Fully Nonlinear})$

2. Homogeneity

Definition 58.7 Homogeneous

$$\mathcal{L}(\mathbf{u}) = 0:$$

All terms depend on \mathbf{u} or on derivatives of \mathbf{u} .

Definition 58.8 Non-Homogeneous

$$\mathcal{L}(\mathbf{u}) = f:$$

Their exists non-zero terms f that do not depend on \mathbf{u} or on derivatives of \mathbf{u} .

3. Constant Coefficients

Definition 58.9 PDEs with Constant Coefficients:

Is a PDE whose coefficients a, b, c, \dots are constants i.e. independent variables.

4. 2nd-Order Linear PDEs in two variables

Definition 58.10

2nd-Order Linear PDEs in two Variables:

$$\mathcal{L}(\mathbf{u}) = a\mathbf{u}_{xx} + 2b\mathbf{u}_{xy} + c\mathbf{u}_{yy} + d\mathbf{u}_x + e\mathbf{u}_y + f = g \quad (58.8)$$

where a, b, \dots, g are functions depending on x and y.

Definition 58.11 Principal Part:

Is the operator \mathcal{L}_0 , that consists of the second-(=highest) order parts of \mathcal{L} :

$$\mathcal{L}_2(\mathbf{u}) := a\mathbf{u}_{xx} + 2b\mathbf{u}_{xy} + c\mathbf{u}_{yy}$$

Definition 58.12 PDEs Discriminante:

Is defined by:

$$\delta(\mathcal{L}) := -\det \begin{pmatrix} a & b \\ b & c \end{pmatrix} = b^2 - ac \quad (58.9)$$

Explanation 58.1.

It turns out that many fundamental properties of the solution of eq. (58.8) are determined by its principal part, or rather by the sign of the discriminant $\delta(\mathcal{L})$.

Definition 58.13

[??]

Parabolic PDEs:

Let [def. 58.10] be a PDE defined on $\Omega \subset \mathbb{R}^2$, then the PDE is called hyperbolic if:

$$\delta(\mathcal{L}) = b^2 - ac = 0 \quad (58.10)$$

Definition 58.14

[??]

Hyperbolic PDEs:

Let [def. 58.10] be a PDE defined on $\Omega \subset \mathbb{R}^2$, then the PDE is called hyperbolic if:

$$\delta(\mathcal{L}) = b^2 - ac > 0 \quad (58.11)$$

Definition 58.15

[??]

Parabolic PDEs:

Let [def. 58.10] be a PDE defined on $\Omega \subset \mathbb{R}^2$, then the PDE is called elliptic if:

$$\delta(\mathcal{L}) = b^2 - ac < 0 \quad (58.12)$$

Explanation 58.2.

The reason for this categorization are normal quadratic equations in two variables:

$$Ax^2 + By^2 + Cxy + Dx + Ey + f = 0$$

If $B^2 - 4AC = 0 \Rightarrow$ the equation is a parabola.

If $B^2 - 4AC > 0 \Rightarrow$ the equation is a hyperbola.

If $B^2 - 4AC < 0 \Rightarrow$ the equation is an ellipse.

2. Method Of Characteristics

Is a method that makes use of geometrical aspects in order to solve 1st-order PDEs with two variables by constructing integral surfaces and can be used to solve PDEs of the type:

$$\text{Linear: } a(x, y)\mathbf{u}_x + b(x, y)\mathbf{u}_y = c(x, y) \quad (58.13)$$

$$\text{Semilin.: } a(x, y)\mathbf{u}_x + b(x, y)\mathbf{u}_y = c(x, y, \mathbf{u}) \quad (58.14)$$

$$\text{Quasilin.: } a(x, y, \mathbf{u})\mathbf{u}_x + b(x, y, \mathbf{u})\mathbf{u}_y = c(x, y, \mathbf{u}) \quad (58.15)$$

Formula 58.1 Method of Characteristics:

$$x := x(r; s) \quad y := y(r; s) \quad z := u(r; s)$$

$$\text{Parameter.: } \lambda(r; s) := x(r; s)\mathbf{e}_x + y(r; s)\mathbf{e}_y + z(r; s)\mathbf{e}_z$$

$$\frac{\partial \lambda}{\partial r}(r; s) = (a, b, c)$$

$$\mathbf{v} := v(x(r; s), y(r; s), z(r; s))$$

$$\frac{\partial x}{\partial r}(r; s) = \dot{x} = a(\lambda(r; s))$$

$$\frac{\partial y}{\partial r}(r; s) = \dot{y} = b(\lambda(r; s))$$

$$\frac{\partial z}{\partial r}(r; s) = \dot{z} = c(\lambda(r; s))$$

Compact:

$$\dot{x} = a(x, y, u) \quad \dot{y} = b(x, y, u) \quad \dot{u} = c(x, y, u)$$

$$\text{I.C.: } x(0; s) = x_0(s) \quad y(0; s) = y_0(s) \quad u(0; s) = u_0(s)$$

2. Homogeneity

Definition 58.7 Homogeneous

$$\mathcal{L}(\mathbf{u}) = 0:$$

All terms depend on \mathbf{u} or on derivatives of \mathbf{u} .

Definition 58.16 Integral Surface

:)

An function $\phi : \mathbb{R}^3 \rightarrow \mathbb{R}$ is a an integral surface of a vector field $\mathbf{V} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ if ϕ is a surface that has in every point a tangent plane containing a vector $\mathbf{v} = (a, b, c)$ of \mathbf{V} .

Corollary 58.3 PDEs and Integral Surfaces:

The solution of a PDE $\mathbf{u}(x, y)$ can be thought of as an integral surface:

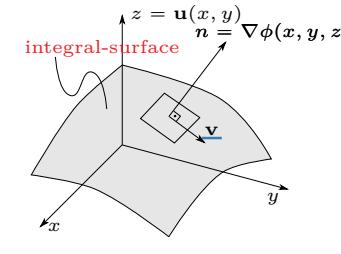
$$z = u(x, y) \quad \text{or implicitly} \quad \phi(x, y, z) = u(x, y) - z \quad (58.16)$$

Explanation 58.3 (

[proof ??]

Integral Surface and PDEs).

The solution $\mathbf{u}(x, y)$ of eq. (58.13) can be sought of as an surface $z = \mathbf{u}(x, y)$ in \mathbb{R}^3 or in implicit form $\phi(x, y, z) := \mathbf{u}(x, y) - z$.



$$\text{Let: } \mathbf{n}(x, y) := \text{grad } \phi = \begin{pmatrix} \phi_x \\ \phi_y \\ \phi_z \end{pmatrix} = \begin{pmatrix} u_x \\ u_y \\ -1 \end{pmatrix} \quad \text{and}$$

$$\text{Let: } \mathbf{V} := \begin{pmatrix} a(x, y) \\ b(x, y) \\ c(x, y) \end{pmatrix} \quad \text{be a vector field } \mathbb{R}^3 \rightarrow \mathbb{R}^3 \text{ and}$$

$$\mathbf{n}(x, y) := \text{grad } \phi = \begin{pmatrix} \phi_x \\ \phi_y \\ \phi_z \end{pmatrix} = \begin{pmatrix} u_x \\ u_y \\ -1 \end{pmatrix}$$

Idea: we can rewrite eq. (58.13) as:

$$\langle (a \ b \ c)^T, \nabla \phi(x, y, z) \rangle = \left\langle \begin{pmatrix} a(x, y) \\ b(x, y) \\ c(x, y) \end{pmatrix}, \begin{pmatrix} u_x \\ u_y \\ -1 \end{pmatrix} \right\rangle = 0$$

Geometric Interpretation:

\mathbf{v} is orthogonal to the normal \mathbf{n} for all points $(x, y, \mathbf{u}(x, y))$.

Hence every vector $\mathbf{v} = (a, b, c)^T$ lies in the tangent plane containing ϕ .

Consequently in order to find a surface ϕ (and thus also a solution \mathbf{u}), we need to search for ϕ s.t. the vector \mathbf{v} lies in the tangent plane for every possible point of ϕ .

Idea

We first simplify the task and start by constructing/finding integral curves λ and then we construct the integral surface ϕ out of this curves.

$$\gamma(r) \quad v_1 \quad v_2 \quad v_3 \quad v_4 \quad v_5$$

3. Linear Equations

Definition 58.17

Characteristic/Integral Curve

$$\lambda_s(r) = \lambda(r; s):$$

Given a vector field \mathbf{V} an integral curve $\lambda(r)$ of that vector field, is a curve parameterized by parameter r :

$$\mathbf{v}(r) := x(r)\mathbf{e}_x + y(r)\mathbf{e}_y + z(r)\mathbf{e}_z = \begin{pmatrix} x(r) \\ y(r) \\ z(r) \end{pmatrix} \quad (58.17)$$

s.t. at each point r of the curve a vector \mathbf{v} of the vector field:

$$\mathbf{v} = \begin{pmatrix} a(x(r), y(r)) \\ b(x(r), y(r)) \\ c(x(r), y(r)) \end{pmatrix} \in \mathbf{V} \quad (58.18)$$

is tangent to the curve:

$$\frac{d\lambda(r)}{dr} = \mathbf{V}(\lambda(r)) = \begin{pmatrix} \frac{dx}{dr} \\ \frac{dy}{dr} \\ \frac{dz}{dr} \end{pmatrix} = \begin{pmatrix} a(\lambda(r)) \\ b(\lambda(r)) \\ c(\lambda(r)) \end{pmatrix} \quad (58.19)$$

Definition 58.18 Characteristic Equations:

The set of ordinary differential equations of a PDE arising from Equation (58.19) are called characteristic equations:

$$\frac{dx(r)}{dr} = \dot{x} = a(\lambda(r)) = a(r) \quad (58.20)$$

$$\frac{dy(r)}{dr} = \dot{y} = b(\lambda(r)) = b(r) \quad (58.21)$$

$$\frac{dz(r)}{dr} = \dot{z} = c(\lambda(r)) = c(r) \quad (58.22)$$

Problem: in order to get a unique solution we need to specify initial conditions.

Idea: If a characteristic has an arbitrary point in common with the integral surface ϕ then the whole characteristic λ will lie in the integral surface.

Proof 58.1: Let: $\phi(\lambda(r)) = u(x(r), y(r)) - z(r)$

$$\Rightarrow \frac{d\phi}{dr} = u_x \frac{dx}{dr} + u_y \frac{dy}{dr} - 1 \frac{dz}{dr} =$$

$$= \begin{pmatrix} u_x \\ u_y \\ -1 \end{pmatrix} \begin{pmatrix} a(x(r), y(r)) \\ b(x(r), y(r)) \\ c(x(r), y(r)) \end{pmatrix} = \begin{pmatrix} u_x \\ u_y \\ -1 \end{pmatrix} \lambda(r) = 0$$

Thus: $\phi(\lambda(r_0)) = 0 \iff \phi(\lambda(r)) = 0, \forall r$

Definition 58.19

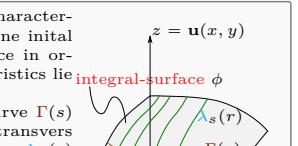
Characteristic (Curve)

$$\lambda_s(r) = \lambda(r; s):$$

is an integral curve of the vector field \mathbf{V} that is uniquely determined by a parameter s .

Consequence: For every characteristic s we need to specify one initial point on the integral surface in order to have all the characteristics lie within the integral surface.

Idea: we define another curve $\Gamma(s)$ on the integral surface that transverses all the characteristic curves $\lambda_s(r)$ transversal (=angle between $\Gamma(s)$ and $\lambda_s(r)$) and $\lambda_s(r)$ is never zero $\Rightarrow \Gamma(s) \not\parallel \lambda_s(r)$.



Definition 58.20 Initial Condition:

$$s \mapsto \Gamma(s), \Gamma : \mathbb{R} \rightarrow \mathbb{R}^3$$

$$\lambda_s(r) = \begin{pmatrix} x_s(r) \\ y_s(r) \\ z_s(r) \end{pmatrix}, \Gamma(s) = \begin{pmatrix} x_0(s) \\ y_0(s) \\ z_0(s) \end{pmatrix} \quad \lambda_s(0) \stackrel{!}{=} \Gamma(s)$$

$$\Rightarrow x_{s_0}(0) = x_0(s) \quad y_{s_0}(0) = y_0(s) \quad z_{s_0}(0) = z_0(s)$$

Definition 58.21 Projected Characteristic Curves

$$\gamma(\tau):$$

Are curves in the plane of the independent variables of our PDE, along which u is constant or satisfies certain conditions. If u is constant along $\gamma(\tau)$ then the initial data is simply propagated along those characteristic curves:

$$\frac{d}{d\gamma} u(\gamma(\tau), \tau) = 0 \iff u(\gamma(\tau), \tau) = u_0(\gamma(\tau)) \quad (58.23)$$

Hint: If the PDE is linear, then the two first characteristics do not depend on u and can be solved directly, u will then be constant along those characteristics:

$$\frac{dx}{dr} = \textcolor{brown}{a} \quad \frac{dy}{dr} = \textcolor{brown}{b} \quad \frac{du}{dr} = \textcolor{brown}{c} \quad \text{implies} \quad \frac{dy}{dx} = \frac{\textcolor{brown}{b}(x, y)}{\textcolor{brown}{a}(x, y)}$$

Hint: If we divide the PDE by $\textcolor{brown}{a}$ we have to solve a PDE less, because the first ODE will always be:

$$\dot{x} = 1 \Rightarrow x = r \Rightarrow \textcolor{teal}{x}_s(r) = x_0(s)$$

4. Quasilinear Equations

Solving Quasilinear Equations

$$\begin{aligned} \textcolor{brown}{a}(x, y, u)\mathbf{u}_x + \textcolor{brown}{b}(x, y, u)\mathbf{u}_y &= \textcolor{brown}{c}(x, y, u) \\ u|_{\Gamma}(r, s) &= \phi(s) \\ \frac{dx}{dr} = \textcolor{brown}{a}(x, y, u) &\quad \frac{dy}{dr} = \textcolor{brown}{b}(x, y, u) \quad \frac{du}{dr} = \textcolor{brown}{c}(x, y, u) \\ \textcolor{teal}{x}_s(0) = x_0(s) &\quad \textcolor{teal}{y}_s(0) = y_0(s) \quad \textcolor{teal}{z}_s(0) = \phi(s) \end{aligned}$$

Results

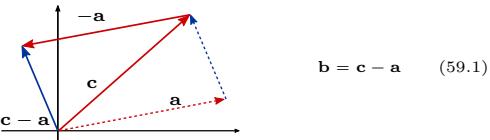
Now the projected characteristic curves may depend on u as well as on x, y . Thus the first two characteristics are no longer decoupled from the third one.

1. We may get projected characteristic curves crossing themselves.
2. u is no longer constant along the projected characteristic curves, rather the PDE reduces to an ODE satisfying certain conditions along these curves.

Linear Algebra

1. Vectors

Definition 59.1 Vector Subtraction:



2. Linear Systems of Equations

1. Gaussian Elimination

2.1.1. Rank

Definition 59.2 Matrix Rank

The ranks of a matrix $A \in \mathbb{R}^{m \times n}$ is defined as the dimension^[def. 59.13] of the vector space spanned^[def. 59.9] by its row or column vectors:

$$\begin{aligned}\text{rank}(A) &= \dim(\{a_1, \dots, a_n\}) \\ &= \dim(\{a_{1,:}, \dots, a_{m,:}\}) \\ \text{def. 59.50} \quad \text{dim}(\mathfrak{R}(A)) &= \dim(\mathfrak{R}(A))\end{aligned}\quad (59.2)$$

Corollary 59.1 :

- The column-and row-ranks of a square matrix $A \in \mathbb{R}^{n \times n}$ are equal.
- The rank of a non-symmetric matrix $A \in \mathbb{R}^{n \times n}$ is limited by the smaller dimension:

$$\text{rank}(A) \leq \min\{n, m\} \quad (59.3)$$

Property 59.1 Rank of Matrix Product:

Let $A \in \mathbb{R}^{m,n}$ and $B \in \mathbb{R}^{n,p}$ then the rank of the matrix product is limited:

$$\text{rank}(AB) \leq \min\{\text{rank}(A), \text{rank}(B)\} \quad (59.4)$$

Rank-1 Matrix

Definition 59.3 Rank-1 Matrix:

Is a matrix of rank one. A tensor product of two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ results in a rank one matrix:

$$\mathbf{uv}^\top = A \in \mathbb{R}^{n,n} \quad (59.5)$$

Definition 59.4 Rank-1 Modification/Update:

Adding a rank-1 matrix to another matrix is called rank-1 modification:

$$X = X + \mathbf{uv}^\top \quad (59.6)$$

3. Sparse Linear Systems

Definition 59.5 Sparse Matrix $A \in \mathbb{K}^{m,n}$, $m, n \in \mathbb{N}_{>0}$:

A matrix A is sparse if:

$$\text{nnz}(A) \ll mn \quad A \in \mathbb{K}^{m,n}, m, n \in \mathbb{N}_{>0} \quad (59.7)$$

$$\text{nnz} := \#\{(i,j) \in \{1, \dots, m\} \times \{1, \dots, n\} : a_{ij} \neq 0\}$$

4. Vector Spaces

1. Vector Space

Definition 59.6 Vector Space: TODO

2. Vector Subspace

Definition 59.7 Vector Subspaces:

A non-empty subset U of a \mathbb{K} -vector space \mathcal{V} is called a subspace of \mathcal{V} if it satisfies:

$$\mathbf{u}, \mathbf{v} \in U \implies \mathbf{u} + \mathbf{v} \in U \quad (59.8)$$

$$\mathbf{u} \in U \implies \lambda \mathbf{u} \in U \quad \forall \lambda \in \mathbb{K} \quad (59.9)$$

Definition 59.8 Linearcombination:

Let $X = \{v_1, \dots, v_n\} \subset \mathcal{V}$ be a non-empty and finite subset of vectors of an \mathbb{K} -vector space \mathcal{V} . A linear combination of X is a combination of the vectors defined as:

$$\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{v}_i = \alpha_1 \mathbf{v}_1 + \dots + \alpha_n \mathbf{v}_n \quad \alpha_i \in \mathbb{K} \quad (59.10)$$

Definition 59.9 Span/Linear Hull

(X) : Is the set of all possible linear combinations^[def. 59.8] of finite set $X = \{v_1, \dots, v_n\} \subset \mathcal{V}$ of a \mathbb{K} vector space \mathcal{V} :

$$(X) = \text{span}(X) = \left\{ \mathbf{v} \mid \sum_{i=1}^n \alpha_i \mathbf{v}_i, \forall \alpha_i \in \mathbb{K} \right\} \quad (59.11)$$

Definition 59.10 Generating Set:

A generating set of vectors $X = \{v_1, \dots, v_m\} \subset \mathcal{V}$ of a vector spaces \mathcal{V} is a set of vectors that span^[def. 59.9] \mathcal{V} :

$$\text{span}(v_1, \dots, v_m) = \mathcal{V} \quad (59.12)$$

Explanation 59.1 (Definition 59.10).

The generating set of vector space (or set of vectors) $\mathcal{V} \stackrel{i.e.}{=} \mathbb{R}^n$ is a subset $X = \{v_1, \dots, v_m\} \subset \mathcal{V}$ s.t. every element of \mathcal{V} can be produced by $\text{span}(X)$.

Definition 59.11 Linear Independence:

A set of vector $\{v_1, \dots, v_n\} \in \mathcal{V}$ is called linear independent if they satisfy:

$$\mathbf{v} = \sum_{i=1}^n \lambda_i \mathbf{v}_i = 0 \iff \alpha_1 = \dots = \alpha_n = 0 \quad (59.13)$$

Corollary 59.2 :

A set of vector $\{x_1, \dots, x_n\} \subset \mathcal{V}$ is called linear independent, if for every subset $X = x_1, \dots, x_m \subseteq \{x_1, \dots, x_n\}$ it holds that:

$$(X) \subseteq \{x_1, \dots, x_n\} \quad (59.14)$$

3. Basis

Definition 59.12 Basis \mathfrak{B} :

A subset $\mathfrak{B} = \{v_1, \dots, v_n\}$ of a \mathbb{K} -vector space \mathcal{V} is called a basis of \mathcal{V} if:

$$\langle \mathfrak{B} \rangle = \mathcal{V} \quad \text{and} \quad \mathfrak{B} \text{ is a linear independent generating set} \quad (59.15)$$

Corollary 59.3 :

The unit vectors e_1, \dots, e_n build a standard basis of the \mathbb{R}^n .

Corollary 59.4 Basis Representation:

Let \mathfrak{B} be a basis of a \mathbb{K} -vector space \mathcal{V} , then it holds that every vector $\mathbf{v} \in \mathcal{V}$ can be represented as a linear combination^[def. 59.8] of \mathfrak{B} by a unique set of coefficients α_i :

$$\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{b}_i \quad \alpha_1, \dots, \alpha_n \in \mathbb{K} \quad b_1, \dots, b_n \in \mathfrak{B} \quad (59.16)$$

4.3.1. Dimensionality

Definition 59.13 Dimension of a vector space $\dim(\mathcal{V})$:

Let \mathcal{V} be a vector space. The dimension of \mathcal{V} is defined as the number of necessary basis vectors $\mathfrak{B} = \{v_1, \dots, v_n\}$ in order to span \mathcal{V} :

$$\dim(\mathcal{V}) := |\mathfrak{B}| = n \in \mathbb{N}_0 \quad (59.17)$$

Corollary 59.5 :

n -linearly independent vectors of a \mathbb{K} -vector space \mathcal{V} with finite dimension n constitute a basis.

Note

If \mathcal{V} is infinite dim ($\mathcal{V}) = \infty$.

4. Affine Subspaces

Definition 59.14 Affine Subspaces:

Given a \mathbb{K} -vector space \mathcal{V} of dimension $\dim(\mathcal{V}) \geq 2$ a sub vector space^[def. 59.7] U of \mathcal{V} defined as:

$$\mathcal{W} := \mathbf{v} + U = \{\mathbf{v} + \mathbf{x} \mid \mathbf{x} \in U\} \quad \mathbf{v} \in \mathcal{V} \quad (59.18)$$

Corollary 59.6 Direction:

The sub vector spaces U are called directions of \mathcal{V} and it holds:

$$\dim(\mathcal{W}) := \dim(U) \quad (59.19)$$

4.4.1. Hyperplanes

Definition 59.15 Hyperplane

\mathcal{H} : A hyperplane is a $d-1$ dimensional subspace of an d -dimensional ambient space that can be specified by the hess normal form^[def. 59.16]:

$$\mathcal{H} = \{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{n}^\top \mathbf{x} - d = 0\} \quad (59.20)$$

Corollary 59.7 Half spaces:

A hyperplane $\mathcal{H} \in \mathbb{R}^{d-1}$ separates its d -dimensional ambient space into two half spaces:

$$\mathcal{H}^+ = \{x \in \mathbb{R}^d \mid \mathbf{n}^\top \mathbf{x} + b > 0\} \quad (59.21)$$

$$\mathcal{H}^- = \{x \in \mathbb{R}^d \mid \mathbf{n}^\top \mathbf{x} + b < 0\} = \mathbb{R}^d - \mathcal{H}^+ \quad (59.22)$$

Notes

Hyperplanes in \mathbb{R}^2 are lines and hyperplanes in \mathbb{R}^3 are lines.

Hess Normal Form

Definition 59.16 Hess Normal Form:

Is an equation to describe hyperplanes^[def. 59.15] in \mathbb{R}^d :

$$\mathbf{r}^\top \mathbf{n} - d = 0 \iff \mathbf{n}^\top (\mathbf{r} - \mathbf{r}_0) \quad \mathbf{r}_0 := \mathbf{r}^\top d \geq 0 \quad (59.23)$$

where all points described by the vector $\mathbf{r} \in \mathbb{R}^d$, that satisfy this equations lie on the hyperplane.

Note

The direction of the unit normal vector is usually chosen s.t. $\mathbf{r}^\top \mathbf{n} \geq 0$.

4.4.2. Lines

Definition 59.17 Lines:

Lines are a set^[def. 50.1] of the form:

$$L = \mathbf{u} + \mathbf{K}\mathbf{v} = \{\mathbf{u} + \lambda \mathbf{v} \mid \lambda \in \mathbb{K}\} \quad \mathbf{u}, \mathbf{v} \in \mathcal{V}, \mathbf{v} \neq 0 \quad (59.24)$$

Two Point Formula

Definition 59.18 Two Point Formula:

$$\begin{array}{c} \text{Diagram showing a line } L \text{ passing through points } u \text{ and } v. \\ \text{The line } L \text{ is defined as } L = \mathbf{u} + \mathbf{K}\mathbf{v}. \end{array} \quad L = \mathbf{u} + \mathbf{K}\mathbf{v} \quad (59.25)$$

4.4.3. Planes

Definition 59.19 Planes:

Planes are sets defined as:

$$E = \mathbf{u} + \mathbf{K}\mathbf{v} + \mathbf{K}\mathbf{w} = \{\mathbf{u} + \lambda \mathbf{v} + \mu \mathbf{w} \mid \lambda, \mu \in \mathbb{K}\} \quad (59.26)$$

$$\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathcal{V} \quad \text{s.t. } \mathbf{v}, \mathbf{u} \neq 0 \quad \text{and} \quad \mathbf{v}, \mathbf{w} \text{ lin. indep.}$$

Parameterform

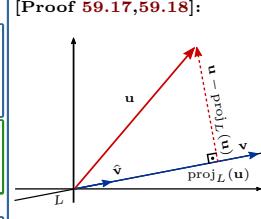
Definition 59.20 Two Point Formula:

$$\begin{array}{c} \text{Diagram showing a plane } E \text{ passing through points } u \text{ and } v. \\ \text{The plane } E \text{ is defined as } E = \mathbf{u} + \mathbf{K}(\mathbf{v} - \mathbf{u}) + \mathbf{K}(\mathbf{w} - \mathbf{u}). \end{array} \quad E = \mathbf{u} + \mathbf{K}(\mathbf{v} - \mathbf{u}) + \mathbf{K}(\mathbf{w} - \mathbf{u}) \quad (59.27)$$

4.4.4. Minimal Distance of Vector Subspaces

Projections in 2D

Definition 59.21 2D Vector Projection [Proof 59.17, 59.18]:



$$\begin{aligned} \mathbf{u}_v &= \text{proj}_L(\mathbf{u}) \\ &= u_v \hat{\mathbf{v}} = (\mathbf{u}^\top \hat{\mathbf{v}}) \hat{\mathbf{v}} \\ &= \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{v}\|^2} \mathbf{v} = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{v}\|^2} \mathbf{v} \end{aligned} \quad (59.28)$$

Proof 59.1: [Corollary 59.8]

$$\frac{1}{\mathbf{v}^\top \mathbf{v}} \mathbf{u}^\top \mathbf{v} \mathbf{v} = \frac{1}{\mathbf{v}^\top \mathbf{v}} \mathbf{v} (\mathbf{v}^\top \mathbf{u}) = \frac{1}{\mathbf{v}^\top \mathbf{v}} \mathbf{v} \mathbf{u} \quad (59.29)$$

General Projections

Definition 59.22 General Vector Projection [proof 59.19]

Is the orthogonal projection \mathbf{u} of a vector \mathbf{v} onto a sub-vector space \mathcal{U}

$$\begin{aligned} \mathbf{u} &= \sum_{i=1}^n \alpha_i \mathbf{b}_i \quad (59.30) \\ \mathbf{AA}^\top \alpha_i &= \mathbf{A}^\top \mathbf{v} \quad \mathbf{A} = \begin{pmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_n \end{pmatrix} \end{aligned}$$

where $\mathfrak{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ is a basis of the vector subspace \mathcal{U} .

Theorem 59.1 Projection Theorem: Let \mathcal{U} a sub vector space of a finite euclidean vector space \mathcal{V} . Then there exists for every vector $\mathbf{v} \in \mathcal{V}$ a vector $\mathbf{u} \in \mathcal{U}$ obtained by an orthogonal projection

$$p : \begin{cases} \mathcal{V} \rightarrow \mathcal{U} \\ \mathbf{v} \mapsto \mathbf{u} \end{cases} \quad (59.31)$$

the vector $\mathbf{u}' := \mathbf{v} - \mathbf{u}$ representing the distance between \mathbf{u} and \mathbf{v} and is minimal:

$$\|\mathbf{u}'\| = \|\mathbf{v} - \mathbf{u}\| \leq \|\mathbf{v} - \mathbf{w}\| \quad \forall \mathbf{w} \in \mathcal{U} \quad \mathbf{u}' \in \mathcal{U}^\perp \quad (59.32)$$

5. Affine Subspaces

6. Planes

<https://math.stackexchange.com/questions/1485509/show-that-two-planes-are-parallel-and-find-the-distance-between-them>

5. Matrices

Special Kind of Matrices

1. Symmetric Matrices

Definition 59.23 Symmetric Matrices: A matrix $\mathbf{A} \in \mathbb{K}^{n \times n}$ is called *symmetric* if it satisfies:

$$\mathbf{A} = \mathbf{A}^T \quad (59.33)$$

Property 59.2 [proof ??] Eigenvalues of real symmetric Matrices: The eigenvalues of a real symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ are real:

$$\text{spectrum}(\mathbf{A}) \in \{\mathbb{R} \geq 0\}_{i=1}^n \quad (59.34)$$

Property 59.3 [proof ??] Orthogonal Eigenvector basis: Eigenvectors of real symmetric matrices with distinct eigenvalues are orthogonal.

Corollary 59.9 Eigendecomposition Symmetric Matrices: If $\mathbf{A} \in \mathbb{R}^{n,n}$ is a real symmetric [def. 59.23] matrix then its eigenvectors are orthogonal and its eigen-decomposition [def. 59.86] is given by:

$$\mathbf{A} = \mathbf{X} \Lambda \mathbf{X}^T \quad (59.35)$$

2. Orthogonal Matrices

Definition 59.24 Orthogonal Matrix: A real valued square matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is said to be orthogonal if its row vectors (and respectively its column vectors) build an orthonormal [def. 59.68] basis:

$$\langle \mathbf{q}_i, \mathbf{q}_j \rangle = \delta_{ij} \quad \text{and} \quad \langle \mathbf{q}_i, \mathbf{q}_j \rangle = \delta_{ij} \quad (59.36)$$

This is exactly true if the inverse of \mathbf{Q} equals its transpose:

$$\mathbf{Q}^{-1} = \mathbf{Q}^T \iff \mathbf{Q}\mathbf{Q}^T = \mathbf{Q}^T\mathbf{Q} = \mathbf{I}_n \quad (59.37)$$

Attention: Orthogonal matrices are sometimes also called orthonormal matrices.

3. Hermitian Matrices

Definition 59.25 Conjugate Transpose $\mathbf{A}^H / \mathbf{A}^*$
Hermitian Conjugate/Adjoint Matrix:

The conjugate transpose of a matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ is defined as:

$$\mathbf{A}^H := (\overline{\mathbf{A}}^T) = \overline{\mathbf{A}}^T \iff a_{i,j}^H = \bar{a}_{j,i} \quad 1 \leq i \leq n \\ 1 \leq j \leq m \quad (59.38)$$

Definition 59.26
Hermitian/Self-Adjoint Matrices $\mathbf{A} = \mathbf{A}^H$:
A hermitian matrix is complex square matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ who is equal to its own conjugate transpose [def. 59.25]:

$$\mathbf{A} = \mathbf{A}^H = \overline{\mathbf{A}}^T \iff a_{i,j} = \bar{a}_{j,i} \quad i \in \{1, \dots, n\} \quad (59.39)$$

Corollary 59.10 : [def. 59.25] implies that \mathbf{A} must be a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$.

Corollary 59.11 Real Hermitian Matrices: From [cor. 50.1] it follows:

$$\mathbf{A} \in \mathbb{R}^{n \times n} \text{ hermitian} \implies \mathbf{A} \text{ real symmetric} \quad (59.40)$$

Property 59.4 [proof 59.15]
Eigenvalues of Hermitian Matrices: The eigenvalues of a hermitian matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ are real:

$$\text{spectrum}(\mathbf{A}) \in \{\mathbb{R} \geq 0\}_{i=1}^n \quad (59.41)$$

Property 59.5 [proof 59.16]
Orthogonal Eigenvector basis: Eigenvectors of hermitian matrices with distinct eigenvalues are orthogonal.

Corollary 59.12 Eigendecomposition Symmetric Matrices: If $\mathbf{A} \in \mathbb{C}^{n,n}$ is a hermitian matrix [def. 59.26] then its eigendecomposition [def. 59.86] is given by:

$$\mathbf{A} = \mathbf{X} \Lambda \mathbf{X}^H \quad (59.42)$$

4. Unitary Matrices

Definition 59.27 Unitary Matrix $\mathbf{U} \mathbf{U}^H$:
is a complex square matrix $\mathbf{U} \in \mathbb{C}^{n \times n}$ whose inverse [def. 59.41] is equal to its conjugate transpose [def. 59.25]:

$$\mathbf{U}^H \mathbf{U} = \mathbf{U} \mathbf{U}^H = \mathbf{I} \quad (59.43)$$

Corollary 59.13 Real Unitary Matrix: A real matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ that is unitary is an orthogonal matrix [def. 59.24].

Property 59.6 Preservation of Euclidean Norm [proof 59.14]: Orthogonal and unitary matrices $\mathbf{Q} \in \mathbb{K}^{n,n}$ do not affect the 2-norm:

$$\|\mathbf{Q}\mathbf{x}\|_2 = \|\mathbf{x}\|_2 \quad \forall \mathbf{x} \in \mathbb{K}^n \quad (59.44)$$

5. Similar Matrices

Definition 59.28 Similar Matrices: Two square matrices $\mathbf{A} \in \mathbb{K}^{n \times n}$ and $\mathbf{B} \in \mathbb{K}^{n \times n}$ are called *similar* if there exists a invertible matrix $\mathbf{S} \in \mathbb{K}^{n \times n}$ s.t.:

$$\exists \mathbf{S}: \quad \mathbf{B} = \mathbf{S}^{-1} \mathbf{A} \mathbf{S} \quad (59.45)$$

Corollary 59.14 Similarity Transformation/Conjugation:
The mapping:

$$\mathbf{A} \mapsto \mathbf{S}^{-1} \mathbf{A} \mathbf{S} \quad (59.46)$$

is called *similarity transformation*

Corollary 59.15 Eigenvalues of Similar Matrices [proof 59.13]:

If $\mathbf{A} \in \mathbb{K}^{n \times n}$ has the eigenvalue-eigenvector pairs $\{\langle \lambda_i, \mathbf{v}_i \rangle\}_{i=1}^n$ then its conjugateeq. (59.46) \mathbf{B} has the same eigenvalues with transformed eigenvectors:

$$\{\langle \lambda_i, \mathbf{u}_i \rangle\}_{i=1}^n \quad \mathbf{u}_i := \mathbf{S}^{-1} \mathbf{v}_i \quad (59.47)$$

6. Block Partitioned Matrices

Definition 59.29 **Skey Symmetric/Antisymmetric Matrices:**

$$\mathbf{A}^T = -\mathbf{A} \quad (59.48)$$

7. Triangular Matrix

Definition 59.30 Triangular Matrix: An upper (lower) triangular matrix, is a matrix whose element's below (above) the main diagonal are all zero:

$$\begin{pmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & & \ddots & \\ l_{n1} & l_{n2} & \cdots & l_{nn} \end{pmatrix} \quad \begin{pmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ u_{21} & u_{22} & \cdots & u_{2n} \\ & & \ddots & \\ & & & u_{nn} \end{pmatrix}$$

Figure 15: Lower Tri. Mat. Figure 16: Upper Tri. Mat.

5.7.1. Unitriangular Matrix

Definition 59.31 Unitriangular Matrix: An upper (lower) unitriangular matrix, is a upper (lower) triangular matrix [def. 59.30] whose diagonal elements are all ones.

5.7.2. Strictly Triangular Matrix

Definition 59.32 Strictly Triangular Matrix: An upper (lower) strictly triangular matrix, is a upper (lower) triangular matrix [def. 59.30] whose diagonal elements are all zero.

8. Block Partitioned Matrices

Definition 59.33 Block Partitioned Matrix: A matrix $\mathbf{M} \in \mathbb{R}^{k+l, k+l}$ can be partitioned into a block partitioned matrix:

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \quad \mathbf{A} \in \mathbb{R}^{k,k}, \mathbf{B} \in \mathbb{R}^{k,l}, \mathbf{C} \in \mathbb{R}^{l,k}, \mathbf{D} \in \mathbb{R}^{l,l} \quad (59.49)$$

Definition 59.34 Block Partitioned Linear System: A linear system $\mathbf{M}\mathbf{x} = \mathbf{b}$ with $\mathbf{M} \in \mathbb{R}^{k+l, k+l}$ and $\mathbf{x}, \mathbf{b} \in \mathbb{R}^{k+l}$ can be partitioned into a block partitioned system:

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} \quad \mathbf{A} \in \mathbb{R}^{k,k}, \mathbf{B} \in \mathbb{R}^{k,l}, \mathbf{C} \in \mathbb{R}^{l,k}, \mathbf{D} \in \mathbb{R}^{l,l} \\ \mathbf{x}_1, \mathbf{b}_1 \in \mathbb{R}^k, \mathbf{x}_2, \mathbf{b}_2 \in \mathbb{R}^l \quad (59.50)$$

5.8.1. Schur Complement

Definition 59.35 Schur Complement: Given a block partitioned matrix [def. 59.33] $\mathbf{M} \in \mathbb{R}^{k+l, k+l}$ its Schur complements are given by:

$$\mathbf{S}_A = \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \quad \mathbf{S}_D = \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} \quad (59.51)$$

5.8.2. Inverse of Block Partitioned Matrix

Definition 59.36 **Inverse of a Block Partitioned Matrix:** Given a block partitioned matrix [def. 59.33] $\mathbf{M} \in \mathbb{R}^{k+l, k+l}$ its inverse \mathbf{M}^{-1} can be partitioned as well:

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \quad \mathbf{M}^{-1} = \begin{bmatrix} \tilde{\mathbf{A}} & \tilde{\mathbf{B}} \\ \tilde{\mathbf{C}} & \tilde{\mathbf{D}} \end{bmatrix} \quad (59.52)$$

$$\begin{aligned} \tilde{\mathbf{A}} &= \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}\mathbf{S}_A^{-1}\mathbf{C}\mathbf{A}^{-1} & \tilde{\mathbf{C}} &= -\mathbf{S}_A^{-1}\mathbf{C}\mathbf{A}^{-1} \\ \tilde{\mathbf{B}} &= -\mathbf{A}^{-1}\mathbf{B}\mathbf{S}_A^{-1} & \tilde{\mathbf{D}} &= \mathbf{S}_A^{-1} \end{aligned}$$

where $\mathbf{S}_A = \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}$ is the Schur complement of \mathbf{A} .

9. Properties of Matrices

9.1. Square Root of p.s.d. Matrices

9.2. Trace

Definition 59.37 Square Root: The trace of an $\mathbf{A} \in \mathbb{R}^{n \times n}$ matrix is defined as:

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii} = a_{11} + a_{22} + \cdots + a_{nn} \quad (59.53)$$

Property 59.7 Trace of a Scalar:

$$\text{tr}(\mathbb{R}) = \mathbb{R} \quad (59.54)$$

Property 59.8 Trace of Transpose:

$$\text{tr}(\mathbf{A}^T) = \text{tr}(\mathbf{A}) \quad (59.55)$$

Property 59.9 Trace of multiple Matrices:

$$\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CBA}) \quad (59.56)$$

6. Matrices and Determinants

1. Determinants

6.1.1. Laplace/Cofactor Expansion

Definition 59.39 Minor:

Definition 59.40 Cofactors:

Properties

Property 59.10 Determinant times Scalar $\det(\alpha \mathbf{A})$: Given a matirx $\mathbf{A} \in \mathbb{R}^{n \times n}$ it holds:

$$\det(\alpha \cdot \mathbf{A}) = \alpha^n \mathbf{A} \quad (59.57)$$

2. Inverse of Matrices

Definition 59.41 Inverse Matrix

$$\mathbf{A}^{-1}$$

6.2.1. Invertability

Definition 59.42

Singular/Non-Invertible Matrix

$$\det(\mathbf{A}) = 0$$

A square matrix $\mathbf{A} \in \mathbb{K}^{n \times n}$ is singular or non-invertible if it satisfies the following and equal conditions:

- $\det(\mathbf{A}) = 0$
- $\mathbf{Ax} = \mathbf{b}$ has either
 - no solution \mathbf{x}
 - infinitely many solutions \mathbf{x}

Transformations And Mapping

7. Linear & Affine Mappings/Transformations

1. Linear Mapping

Definition 59.43

Linear Mapping: A linear mapping, function or transformation is a map $l : V \rightarrow W$ between two \mathbb{K} -vector spaces^[def. 59.6] V and W if it satisfies:

$$l(\mathbf{x} + \mathbf{y}) = l(\mathbf{x}) + l(\mathbf{y}) \quad (\text{Additivity}) \quad (59.58)$$

$$l(\alpha \mathbf{x}) = \alpha l(\mathbf{x}) \quad \forall \alpha \in \mathbb{K} \quad (\text{Homogenity}) \quad (59.59)$$

Proposition 59.1 [proof 59.8]

Equivalent Formulations: Definition 59.43 is equivalent to:

$$l(\alpha \mathbf{x} + \beta \mathbf{y}) = \alpha l(\mathbf{x}) + \beta l(\mathbf{y}) \quad \forall \alpha, \beta \in \mathbb{K} \quad (59.60)$$

Corollary 59.16 Superposition Principle:

Definition 59.43 is also known as the superposition principle: "the net response caused by two or more signals is the sum of the responses that would have been caused by each signal individually."

Corollary 59.17 [proof 59.10]

A linear mapping $\iff \mathbf{Ax}$:

For every matrix $\mathbf{A} \in \mathbb{K}^{m \times n}$ the map:

$$l_{\mathbf{A}} : \begin{cases} \mathbb{K}^n & \rightarrow \mathbb{K}^m \\ \mathbf{x} & \mapsto \mathbf{Ax} \end{cases} \quad (59.61)$$

is a *linear map* and every linear map l can be represented by a matrix vector product:

$$l \text{ is linear} \iff \exists \mathbf{A} \in \mathbb{K}^{n \times m} : f(x) = \mathbf{Ax} \quad \forall x \in \mathbb{K}^m \quad (59.62)$$

Principle 59.1 [proof 59.9]

Principle of linear continuation: A linear mapping $l : \mathcal{V} \rightarrow \mathcal{W}$ is determined by the image of the basis \mathcal{B} of \mathcal{V} :

$$l(\mathbf{v}) = \sum_{i=1}^n \beta_i l(b_i) \quad \mathcal{B}(\mathcal{V}) = \{b_1, \dots, b_n\} \quad (59.63)$$

Property 59.11 [proof 59.11]

Compositions of linear mappings are linear $f \circ g$: Let g, f be linear functions mapping from \mathcal{V} to \mathcal{W} (i.e. matching) then it holds that $f \circ g$ is a linear^[def. 59.43].

2. Affine Mapping

Definition 59.45 Affine Transformation/Map:

Let $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ then:

$$\mathbf{Y} = \mathbf{Ax} + \mathbf{b} \quad (59.64)$$

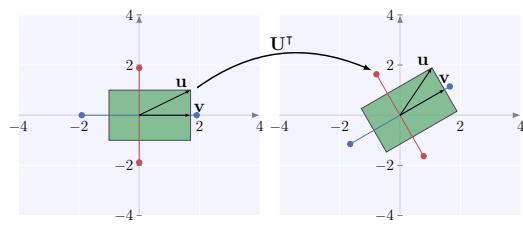
is called an affine transformation of \mathbf{x} .

3. Orthogonal Transformations

Definition 59.46 Orthogonal Transformation:

A linear transformation $T : \mathcal{V} \rightarrow \mathcal{V}$ of an inner product space^[def. 59.78] is an orthogonal transformation if preserves the inner product:

$$T(\mathbf{u}) \cdot T(\mathbf{v}) = \mathbf{u} \cdot \mathbf{v} \quad \forall \mathbf{u}, \mathbf{v} \in \mathcal{V} \quad (59.65)$$



Corollary 59.18 Orthogonal Matrix Transformation: An orthogonal matrix^[def. 59.24] \mathbf{Q} provides an orthogonal transformation:

$$(\mathbf{Qu})^T (\mathbf{Qv}) = \mathbf{uv} \quad (59.66)$$

Explanation 59.2 (Improper Rotations).

Orthogonal transformations in two or three dimensional euclidean space^[def. 59.46] represent improper rotations:

- Stiff Rotations
- Reflections+Rotations
- Reflections

Corollary 59.19 Preservation of Orthogonality: Orthogonal transformation preserver orthogonality.

Corollary 59.20 [proof 59.6]

Preservation of Norm:

An orthogonal transformation $\mathbf{Q} : \mathcal{V} \rightarrow \mathcal{V}$ preserves the length/norm:

$$\|\mathbf{u}\|_{\mathcal{V}} = \|\mathbf{Qu}\|_{\mathcal{V}} \quad (59.67)$$

Corollary 59.21 Preservation of Angle:

An orthogonal transformation T preserves the angle^[def. 59.66] of its vectors:

$$\angle(\mathbf{u}, \mathbf{v}) = \angle(T(\mathbf{u}), T(\mathbf{v})) \quad (59.68)$$

4. Kernel & Image

7.4.1. Kernel

Definition 59.47 Kernel/Null Space $\mathbb{N}/\varphi^{-1}(\{0\})$:

Let φ be a linear mapping^[def. 59.43] between two \mathbb{K} -vector spaces $\varphi : \mathcal{V} \rightarrow \mathcal{W}$.

The *kernel* of φ is defined as:

$$\mathbb{N}(\varphi) := \varphi^{-1}(\{0\}) = \{\mathbf{v} \in \mathcal{V} \mid \varphi(\mathbf{v}) = \mathbf{0}\} \subseteq \mathcal{V} \quad (59.69)$$

Definition 59.48 Right Null Space $\mathbb{N}(\mathbf{A})$:

If $\varphi = \mathbf{A} \in \mathbb{K}^{m \times n}$ then the eq. (59.69) is equal to:

$$\mathbb{N}(\mathbf{A}) = \varphi_{\mathbf{A}}^{-1}(\{0\}) = \{\mathbf{v} \in \mathbb{K}^n \mid \mathbf{Av} = \mathbf{0}\} \in \mathbb{K}^m \quad (59.70)$$

Definition 59.49 Left Null Space $\mathbb{N}(\mathbf{A}^T)$:

If $\varphi = \mathbf{A} \in \mathbb{K}^{m \times n}$ then the left null space is defined as:

$$\mathbb{N}(\mathbf{A}^T) = \varphi_{\mathbf{A}^T}^{-1}(\{0\}) = \{\mathbf{v} \in \mathbb{K}^m \mid \mathbf{A}^T \mathbf{v} = \mathbf{0}\} \in \mathbb{K}^n \quad (59.71)$$

Note

The term *left* null space stems from the fact that:

$$(\mathbf{A}^T \mathbf{x})^T = \mathbf{0} \quad \text{is equal to} \quad \mathbf{x}^T \mathbf{A} = \mathbf{0}$$

7.4.2. Image

Definition 59.50 Image/Range \mathcal{R}/φ :

Let φ be a linear mapping^[def. 59.43] between two \mathbb{K} -vector spaces $\varphi : \mathcal{V} \rightarrow \mathcal{W}$.

The *image* of φ is defined as:

$$\mathcal{R}(\varphi) := \varphi(\mathcal{V}) = \{\varphi(\mathbf{v}) \mid \mathbf{v} \in \mathcal{V}\} \subseteq \mathcal{W} \quad (59.72)$$

Definition 59.51 Column Space \mathbf{Ax} :

If $\varphi = \mathbf{A} = (\mathbf{c}_1, \dots, \mathbf{c}_n) \in \mathbb{K}^{m \times n}$ then eq. (59.72) is equal to:

$$\begin{aligned} \mathcal{R}(\mathbf{A}) &= \varphi_{\mathbf{A}}(\mathbb{K}^n) = \{\mathbf{Ax} \mid \forall \mathbf{x} \in \mathbb{K}^n\} = \langle (\mathbf{c}_1, \dots, \mathbf{c}_n) \rangle \\ &= \left\{ \mathbf{v} \mid \sum_{i=1}^n \alpha_i \mathbf{c}_i, \forall \alpha_i \in \mathbb{K} \right\} \end{aligned} \quad (59.73)$$

Definition 59.52 Row Space $\mathbf{A}^T \mathbf{x}$:

If $\varphi = \mathbf{A} = (r_1^T, \dots, r_m^T) \in \mathbb{K}^{m \times n}$ then the column space is defined as:

$$\begin{aligned} \mathcal{R}(\mathbf{A}^T) &= \varphi_{\mathbf{A}}(\mathbb{K}^m) = \{\mathbf{A}^T \mathbf{x} \mid \forall \mathbf{x} \in \mathbb{K}^m\} = \langle (r_1^T, \dots, r_m^T) \rangle \\ &= \left\{ \mathbf{v} \mid \sum_{i=1}^m \alpha_i r_i^T, \forall \alpha_i \in \mathbb{K} \right\} \end{aligned} \quad (59.74)$$

From orthogonality it follows $x \in \mathcal{R}(\mathbf{A})$, $y \in \mathcal{N}(\mathbf{A}) \Rightarrow x^T y = 0$.

8. Eigenvalues and Vectors

Definition 59.53 Eigenvalues:

Given a square matrix $\mathbf{A} \in \mathbb{K}^{n \times n}$ the eigenvalues

Definition 59.54 Spectrum: The spectrum of a square matrix $\mathbf{A} \in \mathbb{K}^{n \times n}$ is the set of its eigenvalues^[def. 59.53]: spectrum(\mathbf{A}) = $\lambda(\mathbf{A}) = \{\lambda_1, \dots, \lambda_n\}$ (59.78)

Formula 59.1 Eigenvalues of a 2x2 matrix: Given a 2x2-matrix \mathbf{A} its eigenvalues can be calculated by:

$$\{\lambda_1, \lambda_2\} \in \frac{\text{tr}(\mathbf{A}) \pm \sqrt{\text{tr}(\mathbf{A})^2 - 4 \det(\mathbf{A})}}{2} \quad (59.79)$$

with $\text{tr}(\mathbf{A}) = a + d$ $\det(\mathbf{A}) = ad - bc$

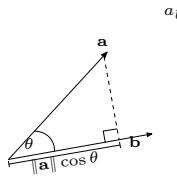
9. Vector Algebra

1. Dot/Standard Scalar Product

Definition 59.55 Scalar Projection

The scalar projection of a vector \mathbf{a} onto a vector \mathbf{b} is the scalar magnitude of the shadow/projection of the vector \mathbf{a} onto \mathbf{b} :

$$a_b = \|\mathbf{a}\| \cos \theta_{\mathbf{a}, \mathbf{b}} = \tilde{\mathbf{a}} \cdot \mathbf{b} \quad (59.80)$$



Definition 59.56 [proof 59.4] **Standard Scalar/Dot Product:**

Given two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{K}^n$ the standard scalar product is defined as:

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \mathbf{v} = \langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^n u_i v_i = u_1 v_1 + \dots + u_n v_n \\ = \|\mathbf{a}\| \|\mathbf{b}\| \cos \theta = u \hat{\mathbf{v}} = v \hat{\mathbf{u}} \quad \theta \in [0, \pi] \quad (59.81)$$

Explanation 59.3 (Geometric Interpretation).

It is the magnitude of one vector times the magnitude of the shadow/scalar projection of the other vector.

Thus the dot product tells you:

1. How much are two vectors pointing into the same direction
2. With what magnitude

Property 59.12 Orthogonal Direction \perp :

For $\theta \in [-\pi, \pi/2]$ $\cos \theta = 0$ and it follows:

$$\mathbf{u} \cdot \mathbf{v} = 0 \iff \mathbf{u} \perp \mathbf{v} \quad (59.82)$$

Note: Perpendicular

Perpendicular corresponds to orthogonality of two lines.

Property 59.13 Maximizing Direction:

For $\theta = 0$ $\cos \theta = 1$ and it follows:

$$\mathbf{u} \cdot \mathbf{v} = \|\mathbf{u}\| \|\mathbf{v}\| \quad (59.83)$$

Property 59.14 Minimizing Direction:

For $\theta = \pi$ $\cos \theta = -1$ and it follows:

$$\mathbf{u} \cdot \mathbf{v} = -\|\mathbf{u}\| \|\mathbf{v}\| \quad (59.84)$$

Definition 59.57 Vector Projection:

2. Cross Product
3. Outer Product

Definition 59.58 Outer Product $\mathbf{u} \mathbf{v}^T = \mathbf{u} \otimes \mathbf{v}$:

Given two vectors $\mathbf{u} \in \mathbb{K}^m$, $\mathbf{v} \in \mathbb{K}^n$ their outer product is defined as:

$$\mathbf{u} \otimes \mathbf{v} = \mathbf{u} \mathbf{v}^H = [\mathbf{u}_1 \dots \mathbf{u}_m] \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_n \end{bmatrix} \quad (59.85) \\ = \begin{bmatrix} \mathbf{u}_1 \odot \mathbf{v}_1 \\ \vdots \\ \mathbf{u}_m \odot \mathbf{v}_n \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 \mathbf{v}_1 & \mathbf{u}_1 \mathbf{v}_2 \dots & \mathbf{u}_1 \mathbf{v}_n \\ \mathbf{u}_2 \mathbf{v}_1 & \mathbf{u}_2 \mathbf{v}_2 \dots & \mathbf{u}_2 \mathbf{v}_n \\ \vdots & \vdots & \vdots \\ \mathbf{u}_m \mathbf{v}_1 & \mathbf{u}_m \mathbf{v}_2 \dots & \mathbf{u}_m \mathbf{v}_n \end{bmatrix}$$

Proposition 59.2 [proof 59.5]

Rank of Outer Product: The outer product of two vectors is of rank one:

$$\text{rank}(\mathbf{u} \otimes \mathbf{v}) = 1 \quad (59.86)$$

4. Vector Norms

Definition 59.59 Norm $\|\cdot\|_{\mathcal{V}}$:

Let \mathcal{V} be a vector space over a field F , a norm on \mathcal{V} is a map:

$$\|\cdot\|_{\mathcal{V}} : \mathcal{V} \mapsto \mathbb{R}_+ \quad (59.87)$$

that satisfies:

$$\|\mathbf{x}\|_{\mathcal{V}} = 0 \iff \mathbf{x} = 0 \quad (\text{Definiteness}) \quad (59.88)$$

$$\|\alpha \mathbf{x}\|_{\mathcal{V}} = |\alpha| \|\mathbf{x}\|_{\mathcal{V}} \quad (\text{Homogeneity}) \quad (59.89)$$

$$\|\mathbf{x} + \mathbf{y}\|_{\mathcal{V}} \leq \|\mathbf{x}\|_{\mathcal{V}} + \|\mathbf{y}\|_{\mathcal{V}} \quad (\text{Triangular Inequality}) \quad (59.90)$$

$$\alpha \in \mathbb{K}$$

$$\forall \mathbf{x}, \mathbf{y} \in \mathcal{V}$$

Explanation 59.4 (Definition 59.59).

A norm is a measure of the size of its argument.

Corollary 59.24 Normed vector space: Is a vector space \mathcal{V} over a field F , on which a norm $\|\cdot\|_{\mathcal{V}}$ can be defined.

9.4.1. Cauchy Schwartz

Definition 59.60 [proof 59.21] **Cauchy Schwartz Inequality:**

$$|\mathbf{u}^T \mathbf{v}| \leq \|\mathbf{u}\| \|\mathbf{v}\| \quad (59.91)$$

9.4.2. Triangular Inequality

Definition 59.61 [proof 59.22] **Triangular Inequality:** States that the length of the sum of two vectors is lower or equal to the sum of their individual lengths:

$$\|\mathbf{u} + \mathbf{v}\| \leq \|\mathbf{u}\| + \|\mathbf{v}\| \quad (59.92)$$

Corollary 59.25 Reverse Triangular Inequality:

$$-\|\mathbf{x} - \mathbf{y}\|_{\mathcal{V}} \leq \|\mathbf{x}\|_{\mathcal{V}} - \|\mathbf{y}\|_{\mathcal{V}} \leq \|\mathbf{x} - \mathbf{y}\|_{\mathcal{V}} \\ \text{resp. } \|\mathbf{x}\|_{\mathcal{V}} - \|\mathbf{y}\|_{\mathcal{V}} \leq \|\mathbf{x} - \mathbf{y}\|_{\mathcal{V}}$$

5. Distances

Definition 59.62

Distance Function/Measure $d : S \times S \mapsto \mathbb{R}_+$: Let S be a set, a distance function is a mapping d that satisfies:

$$d(\mathbf{x}, \mathbf{x}) = 0 \quad (\text{Zero Identity Distance}) \quad (59.93)$$

$$d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}) \quad (\text{Symmetry}) \quad (59.94)$$

$$d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}) \quad (\text{Triangular Identity}) \quad (59.95)$$

$$\forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in S$$

Explanation 59.5 (Definition 59.62).

Is measuring the distance between two things.

9.5.1. Contraction

Definition 59.63 **Contraction:** Given a metric space (M, d) is a mapping $f : M \mapsto M$ that satisfies:

$$d(f(\mathbf{x}), f(\mathbf{y})) \leq \lambda d(\mathbf{x}, \mathbf{y}) \quad \lambda \in [0, 1] \quad (59.96)$$

6. Metrics

Definition 59.64 **Metric** $d : S \times S \mapsto \mathbb{R}_+$:

Is a distance measure [def. 59.62] that additionally satisfies the identity of indiscernibles:

$$d(\mathbf{x}, \mathbf{y}) = 0 \iff \mathbf{x} = \mathbf{y} \quad \forall \mathbf{x}, \mathbf{y} \in S$$

Corollary 59.26 Metric→Norm: Every norm $\|\cdot\|_{\mathcal{V}}$ on a vector space \mathcal{V} over a field F induces a metric by:

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_{\mathcal{V}} \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{V}$$

metric induced by norms additionally satisfy: $\forall \mathbf{x}, \mathbf{y} \in \mathcal{V}, \alpha \in F \subseteq \mathbb{K}, K = \mathbb{R}$ or \mathbb{C}

1. Homogeneity/Scaling: $d(a\mathbf{x}, a\mathbf{y})_{\mathcal{V}} = |\alpha| d(\mathbf{x}, \mathbf{y})_{\mathcal{V}}$

2. Translational Invariance: $d(\mathbf{x} + \alpha, \mathbf{y} + \alpha) = d(\mathbf{x}, \mathbf{y})$

Conversely not every metric induces a norm but if a metric d on a vector space \mathcal{V} satisfies the properties then it induces a norm of the form:

$$\|\mathbf{x}\|_{\mathcal{V}} := d(\mathbf{x}, 0)_{\mathcal{V}}$$

Note

Similarity measure is a much weaker notion than a metric as triangular inequality does not have to hold.

Hence: If a is similar to b and b is similar to c it does not imply that a is similar to c .

Note

That similarity measure is a much weaker notion than a metric as triangular inequality does not have to hold.

Hence: If a is similar to b and b is similar to c it does not imply that a is similar to c .

Note

(bilinear form $\xrightarrow{\text{induces}}$)
inner product $\xrightarrow{\text{induces}}$ norm $\xrightarrow{\text{induces}}$ metric.

9.6.1. Metric Space

Definition 59.65 **Metric Space** (M, d) :

A metric space is a pair (M, d) of a set M and a metric d defined on M :

$$d : M \times M \mapsto \mathbb{R}_+ \quad (59.97)$$

10. Angles

Definition 59.66 **Angle between Vectors** $\angle(\mathbf{u}, \mathbf{v})$: Let $\mathbf{u}, \mathbf{v} \in \mathbb{K}^n$ be two vectors of an inner product space [def. 59.78].

\mathcal{V} . The angle $\alpha \in [0, \pi]$ between \mathbf{u}, \mathbf{v} is defined by:

$$\angle(\mathbf{u}, \mathbf{v}) := \alpha \quad \cos \alpha = \frac{\mathbf{u}^T \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \quad \alpha \in [0, \pi] \quad (59.98)$$

11. Orthogonality

Definition 59.67 **Orthogonal Vectors**: Let \mathcal{V} be an inner-product space [def. 59.78]. A set of vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_n\} \subset \mathcal{V}$ is called orthogonal iff:

$$\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0 \quad \forall i \neq j \quad (59.99)$$

1. Orthonormality

Definition 59.68 **Orthonormal Vectors**: Let \mathcal{V} be an inner-product space [def. 59.78]. A set of vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_n\} \subset \mathcal{V}$ is called orthonormal iff:

$$\langle \mathbf{u}_i, \mathbf{u}_j \rangle = \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \quad \forall i, j \quad (59.100)$$

12. Special Kind of Vectors

1. Binary/Boolean Vectors

Definition 59.69 **Binary/Boolean Vectors/Bit Maps** \mathbb{B}^n : Are vectors that contain only zero or one values:

$$\mathbb{B}^n = \{0, 1\}^n \quad (59.101)$$

Definition 59.70 **R-Sparse Boolean Vectors** \mathbb{B}_r^n :

Are boolean vectors that contain exactly r one values:

$$\mathbb{B}_r^n = \left\{ \mathbf{x} \in \{0, 1\}^n : \mathbf{x}^T \mathbf{x} = \sum_{i=1}^n x_i = r \right\} \quad (59.102)$$

2. Probabilistic Vectors

Definition 59.71 **Probabilistic Vectors**: Are vectors that represent probabilities and satisfy:

$$\left\{ \mathbf{x} \in [0, 1]^n : \sum_{i=1}^n x_i = 1 \right\} \quad (59.103)$$

13. Vector Spaces and Measures

1. Bilinear Forms

2. Quadratic Forms

13.2.1. Min/Max Value

Corollary 59.27 [proof 59.20] Extreme Value: The minimum/maximum of a quadratic form?? with a quadratic matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ is given by the eigenvector corresponding to the smallest/largest eigenvalue of \mathbf{A} :

$$\mathbf{v}_1 \in \arg \min_{\mathbf{x}^T \mathbf{x}=1} \mathbf{x}^T \mathbf{A} \mathbf{x} \quad \mathbf{v}_1 \in \arg \max_{\mathbf{x}^T \mathbf{x}=1} \mathbf{x}^T \mathbf{A} \mathbf{x} \quad (59.104)$$

Note

$$(\mathbf{Q}^T \mathbf{n})^T \mathbf{Q}^T \mathbf{n} = \mathbf{n}^T \mathbf{Q} \mathbf{Q}^T \mathbf{n} = \mathbf{n}^T \mathbf{n} = 1$$

13.2.2. Skew Symmetric Matrix

Corollary 59.28

Quadratic Form of Skew Symmetric matrix: The quadratic form of a skew symmetric matrix [def. 59.29] vanishes:

$$\mathbf{a} = \mathbf{x}^T \mathbf{A}_{\text{skew}} \mathbf{x} = (\mathbf{x}^T \mathbf{A}_{\text{skew}}^T \mathbf{x})^T = (\mathbf{x}^T \mathbf{A}_{\text{skew}} \mathbf{x})^T = -\mathbf{a} \quad (59.105)$$

Which can only hold iff $\mathbf{a} = 0$.

3. Inner Product – Generalization of the dot product

Definition 59.72 **Bilinear Form/Functional:**

Is a mapping $a : \mathcal{V} \times \mathcal{V} \mapsto F$ on a field of scalars $F \subseteq \mathbb{K}$, $K = \mathbb{R}$ or \mathbb{C} that satisfies:

$$a(\alpha \mathbf{u} + \beta \mathbf{v}, w) = \alpha a(\mathbf{u}, w) + \beta a(\mathbf{v}, w) \\ a(\mathbf{u}, \alpha \mathbf{v} + \beta \mathbf{w}) = \alpha a(\mathbf{u}, \mathbf{v}) + \beta a(\mathbf{u}, \mathbf{w}) \quad \forall \alpha, \beta \in \mathbb{K}, \mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathcal{V}$$

Thus: a is linear w.r.t. each argument.

Definition 59.73 **Symmetric bilinear form:** A bilinear form a on \mathcal{V} is symmetric if and only if:

$$a(\mathbf{u}, \mathbf{u}) > 0 \quad \forall \mathbf{u} \in \mathcal{V} \setminus \{0\} \quad (59.106)$$

And positive semidefinite $\iff \geq$ (59.107)

Corollary 59.29 **Matrix induced Bilinear Form:**

For finite dimensional inner product spaces $\mathcal{X} \in \mathbb{K}^n$ any symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ induces a bilinear form:

$$a(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' = (\mathbf{A} \mathbf{x}') \mathbf{x}$$

Definition 59.75 **Positive (semi) definite Matrix** \gg :

A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive definite if and only if:

$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \iff \mathbf{A} > 0 \quad \forall \mathbf{x} \in \mathbb{R}^n \setminus \{0\} \quad (59.108)$$

And positive semidefinite $\iff \geq$ (59.109)

Corollary 59.30 [proof 59.2] **Eigenvalues of positive (semi) definite matrix:**

A positive definite matrix is a matrix where every eigenvalue is strictly positive and positive semi definite if every eigenvalue is positive.

$$\forall \lambda_i \in \text{eigen}(\mathbf{A}) > 0 \quad (59.110)$$

And positive semidefinite $\iff \geq$ (59.111)

Note

Positive definite matrices are often assumed to be symmetric but that is not necessarily true.

Proof 59.2: ?? 59.2 (for real matrices):

Let \mathbf{v} be an eigenvector of \mathbf{A} then it follows:

$$\mathbf{v}^T \mathbf{A} \mathbf{v} = \mathbf{v}^T \mathbf{A} \mathbf{v} = \|\mathbf{v}\|^2 \lambda$$

Corollary 59.31 **Positive Definiteness and Determinant:** The determinant of a positive definite matrix is always positive. Thus a positive definite matrix is always nonsingular

lar

Definition 59.76 Negative (semi) definite Matrix <:
A matrix $A \in \mathbb{R}^{n \times n}$ is negative definite if and only if:
 $x^T Ax < 0 \iff A < 0 \quad \forall x \in \mathbb{R}^n \setminus \{0\}$ (59.112)

And negative semidefinite \iff (59.113)

Theorem 59.3 Sylvester's criterion: Let A be symmetric/Hermitian matrix and denote by $A^{(k)}$ the $k \times k$ upper left sub-matrix of A . Then it holds that:

- $A > 0 \iff \det(A^{(k)}) > 0 \quad k = 1, \dots, n$ (59.114)
- $A < 0 \iff (-1)^k \det(A^{(k)}) > 0 \quad k = 1, \dots, n$ (59.115)

- A is indefinite if the first $\det(A^{(k)})$ that breaks both of the previous patterns is on the wrong side.
- Sylvester's criterion is inconclusive (A can be anything of the previous three) if the first $\det(A^{(k)})$ that breaks both patterns is 0.

14. Inner Products

Definition 59.77 Inner Product: Let \mathcal{V} be a vector space over a field $F \subseteq \mathbb{K}$ of scalars. An inner product on \mathcal{V} is a map:
 $\langle \cdot, \cdot \rangle : \mathcal{V} \times \mathcal{V} \mapsto F \subseteq \mathbb{K} \quad K = \mathbb{R} \text{ or } \mathbb{C}$ (59.116)

that satisfies: $\forall x, y, z \in \mathcal{V}, \alpha, \beta \in F$

1. (Conjugate) Symmetry: $\langle x, y \rangle = \langle y, x \rangle$.

2. Linearity in the first argument:

$$\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$$

3. Positive-definiteness:

$$\langle x, x \rangle \geq 0 : x = 0 \iff \langle x, x \rangle = 0$$

Definition 59.78 Inner Product Space $(\mathcal{V}, \langle \cdot, \cdot \rangle_{\mathcal{V}})$: Let $F \subseteq \mathbb{K}$ be a field of scalars.

An inner product space \mathcal{V} is a vector space over a field F together with an an inner product $\langle \cdot, \cdot \rangle_{\mathcal{V}}$.

Corollary 59.32 Inner product \rightarrow S.p.d. Bilinear Form: Let \mathcal{V} be a vector space over a field $F \subseteq \mathbb{K}$ of scalar.

An inner product on \mathcal{V} is a positive definite symmetric bilinear form on \mathcal{V} .

Example: scalar prodct

Let $a(u, v) = u^T Iv$ then the standard scalar product can be defined in terms of a bilinear form vice versa the standard scalar product induces a bilinear form.

Note

Inner products must be positive definite by defintion $\langle x, x \rangle \geq 0$, whereas bilinear forms must not.

Corollary 59.33 Inner product induced norm $\langle \cdot, \cdot \rangle_{\mathcal{V}} \rightarrow \|\cdot\|_{\mathcal{V}}$: Every inner product $\langle \cdot, \cdot \rangle_{\mathcal{V}}$ induces a norm of the form:

$$\|x\|_{\mathcal{V}} = \sqrt{\langle x, x \rangle} \quad x \in \mathcal{V}$$

Thus We can define function spaces by their associated norm $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$ and inner product spaces lead to normed vector spaces and vice versa.

Corollary 59.34 Energy Norm: A s.p.d. bilinear form $a : \mathcal{V} \times \mathcal{V} \mapsto F$ induces an energy norm:

$$\|x\|_a := (a(x, x))^{\frac{1}{2}} = \sqrt{a(x, x)} \quad x \in \mathcal{V}$$

15. Matrix Algebra

16. Matrix Norms

1. Operator Norm

Definition 59.79 Operator/Induced Norm:

Let $\|\cdot\|_{\mu} : \mathbb{K}^m \mapsto \mathbb{R}$ and $\|\cdot\|_{\nu} : \mathbb{K}^n \mapsto \mathbb{R}$ be vector norms. The operator norm is defined as:

$$\|A\|_{\mu, \nu} := \sup_{\substack{x \in \mathbb{K}^n \\ x \neq 0}} \frac{\|Ax\|_{\mu}}{\|x\|_{\nu}} = \sup_{\substack{x \in \mathbb{K}^n \\ x \neq 0}} \frac{\|Ax\|_{\mu}}{\|x\|_{\nu}=1} \quad \| \cdot \|_{\mu} : \mathbb{K}^m \mapsto \mathbb{R}$$

$$(59.117)$$

Explanation 59.6 (Definition 59.79). Is a measure for the largest factor by which a matrix A can stretch a vector $x \in \mathbb{R}^n$.

2. Induced Norms

Corollary 59.35 Induced Norms: Let $\|\cdot\|_p : \mathbb{K}^{m \times n} \mapsto \mathbb{R}$ defined as:

$$\|A\|_p := \sup_{\substack{x \in \mathbb{K}^n \\ x \neq 0}} \frac{\|Ax\|_p}{\|x\|_p} = \sup_{\substack{y \in \mathbb{K}^m \\ y \neq 0}} \frac{\|Ay\|_p}{\|y\|_p=1} \quad (59.118)$$

Explanation 59.7 ([Corollary 59.35]).

Induced norms are matrix norms induced by vector norms as we:

- Only work with vectors Ax
- And use the normal p -vector norms $\|\cdot\|_p$

Note supremum

The set of vectors $\{y \mid \|y\|_p = 1\}$ is compact, thus if we consider finite matrices the supremum is attained and we may replace it by the max.

3. Induced Norms

16.3.1. 1-Norm

Definition 59.80 Column Sum Norm $\|A\|_1 :$

$$\|A\|_1 = \sup_{\substack{x \in \mathbb{K}^n \\ x \neq 0}} \frac{\|Ax\|_1}{\|x\|_1} = \max_{1 \leq i \leq n} \sum_{j=1}^m |a_{ij}| \quad (59.119)$$

16.3.2. ∞ -Norm

Definition 59.81 Row Sum Norm $\|A\|_{\infty} :$

$$\|A\|_{\infty} = \sup_{\substack{x \in \mathbb{K}^n \\ x \neq 0}} \frac{\|Ax\|_{\infty}}{\|x\|_{\infty}} = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \quad (59.120)$$

16.3.3. Spectral Norm

Spectral Radius & Singular Value

Definition 59.82 Spectral Radius $\rho(A) :$

The spectral radius is defined as the largest eigenvalue of a matrix:

$$\rho(A) = \max \{ \lambda \mid \lambda \text{ eigenval}(A) \} \quad (59.121)$$

Definition 59.83 Singular Value $\sigma_i :$

Given a matrix $A \in \mathbb{K}^{m \times n}$ its n real and positive singular values are defined as: $\sigma(A) := \{ \sqrt{\lambda_i} \}_{i=1}^n \mid \lambda_i \in \text{eigenval}(A^T A) \}$ (59.122)

Spectral Norm

Definition 59.84 L2/Spectral Norm $\|A\|_2 :$

$$\|A\|_2 = \sup_{\substack{x \in \mathbb{K}^n \\ \|x\|_2=1}} \|Ax\|_2 = \max_{\|x\|_2=1} \sqrt{x^T A^T A x} \quad (59.123)$$

$$= \max_{\|x\|_2=1} \sqrt{\rho(A^T A)} =: \sigma_{\max}(A) \quad (59.124)$$

4. Energy Norm 5. Forbenius Norm

Definition 59.85 Forbenius Norm $\|A\|_F :$

The Forbenius norm $\|\cdot\|_F : \mathbb{K}^{m \times n} \mapsto \mathbb{R}$ is defined as:

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{tr}(A^H A)} \quad (59.125)$$

6. Distance

17. Decompositions

1. Eigen/Spectral decomposition

Definition 59.86 $A = X \Lambda X^{-1}$, [proof 59.25]

Eigendecomposition/ Spectral Decomposition :

Let $A \in \mathbb{K}^{n \times n}$ be a diagonalizable square matrix and define by $X = [x_1, \dots, x_n] \in \mathbb{K}^{n \times n}$ a non-singular matrix whose column vectors are the eigenvectors of A with associated eigenvalue matrix $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Then A can be represented as:

$$A = X \Lambda X^{-1} \quad (59.126)$$

Proposition 59.3 Diagonalization: If none of A eigenvalues are zero it can be diagonalized:

$$S^{-1} A S = \Lambda \quad (59.127)$$

Proposition 59.4 Existence:

$$\exists X \Lambda X^{-1} \iff A \text{ diagonalizable} \quad (59.128)$$

2. QR-Decompositions

3. Singular Value Decomposition

Definition 59.87

Singular Value Decomposition (SVD) $U \Sigma V^H :$

For any matrix $A \in \mathbb{K}^{m,n}$ there exist unitary matrices^[def. 59.27]

$$U \in \mathbb{K}^{m,m} \quad V \in \mathbb{K}^{n,n}$$

and a (generalized) diagonal matrix:

$$\rho := \min\{m, n\}$$

$$\Sigma = \text{gendiag}(\sigma_1, \dots, \sigma_{\rho}) \in \mathbb{R}^{m,n}$$

such that:

$$A = U \Sigma V^H \quad (59.129)$$

$$= \begin{pmatrix} u_1 & u_2 & u_3 & \dots & u_m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \hline & M_{\text{null}(A)} & & & \end{pmatrix} \begin{pmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_r & & \\ & & & 0 & \\ & & & & 0 \end{pmatrix} \begin{pmatrix} v_1^T & v_2^T & v_3^T & \dots & v_n^T \end{pmatrix}$$

17.3.1. Eigenvalues

Proposition 59.5 [proof 59.23]: The eigenvalues of a matrix $A^T A$ are positive.

Proposition 59.6 [proof 59.24]

Similarity Transformation: The unitary matrix V provides a similarity transformation^[cor. 59.14] of $A^T A$ into a diagonal matrix $\Sigma^T \Sigma$:

$$\Sigma^T \Sigma \mapsto V^H A^T A V \quad (59.130)$$

Corollary 59.36 eigenval($A^T A$) = eigenval($\Sigma^T \Sigma$):

From proposition 59.6 and [cor. 59.15] it follows that:

$$\text{eigenval}(A^T A) = \text{eigenval}(\Sigma^T \Sigma) \quad (59.131)$$

$$\implies \|A\|_2 = \sqrt{\rho(A^T A)} = \sqrt{\lambda_{\max}} = \sigma_{\max}$$

Note

λ and singularvalue corresponds to the eigenvalues/singular-values of $A^T A$ and not A

17.3.2. Best Lower Rank Approximation

Theorem 59.4 Eckart Young Theorem: Given a matrix $X \in \mathbb{K}^{m,n}$ the reduced SVD X defined as:

$$U_k := [u_{:,1} \dots u_{:,k}] \in \mathbb{K}^{m,k}$$

$$X_k := U_k \Sigma_k V_k^H \quad \Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{k,k}$$

$$V_k = [v_{:,1} \dots v_{:,k}] \in \mathbb{K}^{n,k}$$

provides the best lower k rank approximation of X :

$$\min_{Y \in \mathbb{K}^{m,n}, \text{rank}(Y) \leq k} \|X - Y\|_F = \|X - X_k\|_F \quad (59.132)$$

18. Matrix Calculus

1. Derivatives

$$\frac{\partial}{\partial x} (b^T x) = \frac{\partial}{\partial x} (x^T b) = b$$

$$\frac{\partial}{\partial x} (x^T x) = 2x$$

$$\frac{\partial}{\partial x} Ax = A \quad (59.133)$$

$$\frac{\partial}{\partial x} x^T Ax = (A + A^T)x \quad (59.134)$$

$$\frac{\partial}{\partial x} (b^T Ax) = A^T b \quad \frac{\partial}{\partial x} (c^T x) = cb^T \quad \frac{\partial}{\partial x} (\|x - b\|_2) =$$

$$\frac{\partial}{\partial x} (\|x\|_2^2) = \frac{\partial}{\partial x} (x^T x) = 2x \quad \frac{\partial}{\partial x} (\|x\|_F^2) = 2x$$

$$\frac{\partial}{\partial x} (\|Ax - b\|_2^2) = 2(A^T Ax - A^T b) \quad \frac{\partial}{\partial x} (\|x\|) = |x| \cdot x^{-1}$$

19. Proofs

Proof 59.3: [def. 59.36]

$$MM^{-1} = \begin{bmatrix} I_{k,k} & 0_{k,l} \\ 0_{l,k} & I_{l,l} \end{bmatrix} \quad (59.135)$$

1. Vector Algebra

Proof 59.4 Definition 59.56:

$$(1): \|a - b\|^2 \stackrel{\text{eq. (60.19)}}{=} \|a\|^2 + \|b\|^2 - 2\|a\|\|b\| \cos \theta$$

$$(2): \|a - b\| = (a - b)(a - b) = \|a\|^2 + \|b\|^2 - 2(a \cdot b)$$

$$\|a - b\| = \|a - b\| \Rightarrow ab = \|a\|\|b\| \cos \theta$$

Proof 59.5 Proposition 59.2: The outer product of u with v corresponds to a scalar multiplication of v with elements u_i thus the rank must be that of v , which is a vector and hence of rank 1

$$u \otimes v = uv^H = \begin{bmatrix} u_1 \odot v_1 \\ \vdots \\ u_m \odot v_n \end{bmatrix}$$

2. Mappings

Proof 59.6: Corollary 59.20

$$\|Qx\|^2 = (Qx)^T Qx = x^T Q^T Qx = x^T x = \|x\|^2$$

Proof 59.7: Corollary 59.21 Follows immediately from definition 59.66 in combination with eqs. (59.65) and (59.67).

Proof 59.8: Proposition 59.1:

$$\begin{aligned} l(\alpha x + \beta y) &\stackrel{59.58}{=} l(\alpha x) + l(\beta y) \stackrel{59.59}{=} \alpha l(x) + \beta l(y) \\ l(\alpha x + 0) &= \alpha l(x) \\ l(1x + 1y) &= l(x) + l(y) \end{aligned}$$

Proof 59.9 principle 59.1:

Every vector $v \in \mathcal{V}$ can be represented by a basis eq. (59.16) of \mathcal{V} . With homogeneityeq. (59.59) and additivityeq. (59.58) it follows for the image of all $v \in \mathcal{V}$:

$$l(v) = l(\alpha_1 b_1 + \dots + \alpha_n b_n) = l\alpha_1(b_1) + \dots + l(\alpha_n) b_n \quad (59.136)$$

\Rightarrow the image of the basis of \mathcal{V} determines the linear mapping.

Proof 59.10 Proof [Corollary 59.17]:

$$\begin{aligned} \Rightarrow l(A(\alpha x + \beta y)) &= A(\alpha x + \beta y) = \alpha Ax + \beta Ay = \alpha l(x) + \beta l(y) \\ \Leftrightarrow \text{Let } \mathcal{B} \text{ be a standard normal basis of } \mathcal{V} \text{ with eq. (59.136):} \end{aligned}$$

$$l(x) = \sum_{i=1}^n x_i l(e_i) = \sum_{i=1}^n x_i A_{:,i} = Ax \quad A_{:,i} := l(e_i) \in \mathbb{R}^n$$

Proof 59.11 Proof Property 59.11:

$$(g \circ f)(\alpha x) = g(f(\alpha x)) = g(f(\alpha x)) = \alpha(g \circ f)(x)$$

$$(g \circ f)(x+y) = g(f(x+y)) = g(f(x) + f(y))$$

$$= (g \circ f)(x) + (g \circ f)(y)$$

or even simpler as every linear form can be represented by a matrix product:

$$f(y) = Ay \quad g(z) = Bz \quad \Rightarrow \quad (f \circ g)(x) = ABx := Cx$$

Proof 59.12: [Corollary 59.22] Let $\mathbf{y} \in \mathbb{N}(\mathbf{A})$ ($\mathbf{z} \in \mathbb{N}(\mathbf{A}^\top)$) then it follows:

$$\begin{aligned}\mathbb{N}(\mathbf{A}) \perp \mathbb{R}(\mathbf{A}^\top) \quad (\mathbf{A}^\top \mathbf{x})^\top \mathbf{y} = \mathbf{x}^\top \mathbf{A} \mathbf{y} = \mathbf{x}^\top \mathbf{0} = 0 \\ \mathbb{N}(\mathbf{A}^\top) \perp \mathbb{R}(\mathbf{A}) \quad (\mathbf{A} \mathbf{x})^\top \mathbf{z} = \mathbf{x}^\top \mathbf{A}^\top \mathbf{z} = \mathbf{x}^\top \mathbf{0} = 0\end{aligned}$$

3. Special Matrices

Proof 59.13 [Corollary 59.15]: Let $\mathbf{u} = \mathbf{S}^{-1}\mathbf{v}$ then it follows:

$$\mathbf{S}^{-1}\mathbf{A}\mathbf{S}\mathbf{u} = \mathbf{S}^{-1}\mathbf{A}\mathbf{S}\mathbf{v} = \lambda \mathbf{S}^{-1}\mathbf{v} = \lambda \mathbf{u}$$

Proof 59.14 Property 59.6:

$$\|\mathbf{Q}\mathbf{x}\|_2^2 = (\mathbf{Q}\mathbf{x})^\top \mathbf{Q}\mathbf{x} = \mathbf{x}^\top \mathbf{Q}^\top \mathbf{Q}\mathbf{x} = \|\mathbf{x}\|_2^2$$

Proof 59.15: Property 59.4

Let $\mathbf{A} \in \mathbb{K}^{n \times n}$ be a hermitian matrix^[def. 59.26] and let $\lambda \in \mathbb{K}$ be an eigenvalue of \mathbf{A} with corresponding eigenvector $\mathbf{v} \in \mathbb{K}^n$:

$$\begin{aligned}\lambda(\bar{\mathbf{v}}^\top \mathbf{v}) &= \bar{\mathbf{v}}^\top \lambda \mathbf{v} = \bar{\mathbf{v}}^\top \mathbf{A}\mathbf{v} = (\bar{\mathbf{v}}^\top \mathbf{A}\mathbf{v})^\top = \bar{\mathbf{A}}\bar{\mathbf{v}}^\top \mathbf{v} = \bar{\lambda}(\bar{\mathbf{v}}^\top \mathbf{v}) \\ \lambda(\bar{\mathbf{v}}^\top \mathbf{v}) &= \bar{\lambda}(\bar{\mathbf{v}}^\top \mathbf{v})\end{aligned}$$

$$1. \bar{\mathbf{v}}^\top \mathbf{v} = \sum_{i=1}^n |v_i|^2 > 0 \text{ as } \mathbf{v} \neq \mathbf{0}$$

$$2. \lambda = \bar{\lambda} \text{ which can only hold for } \lambda \in \mathbb{R} \text{ (Equation (50.8))}$$

Proof 59.16: ??

4. Vector Spaces

Proof 59.17 Definition 59.21: We know that $\text{proj}_L(\mathbf{u})$ must be a vector times a certain magnitude:

$$\text{proj}_L(\mathbf{u}) = \alpha \bar{\mathbf{v}} \quad \alpha \in \mathbb{K} \quad (59.137)$$

the magnitude follows from the scalar projection^[def. 59.55] in the direction of \mathbf{v} which concludes the derivation.

Proof 59.18 Definition 59.21 (via orthogonality): We know that $\mathbf{u} - \text{proj}_L(\mathbf{u})$ must be orthogonal^[def. 59.67] to \mathbf{v}

$$(\mathbf{u} - \text{proj}_L(\mathbf{u}))^\top \mathbf{v} = (\mathbf{u} - \alpha \bar{\mathbf{v}})^\top \mathbf{v} = 0 \Rightarrow \alpha = \frac{\mathbf{u}^\top \mathbf{v}}{\mathbf{v}^\top \mathbf{v}}$$

Proof 59.19: Definition 59.22 Let $\mathcal{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ a basis of \mathcal{U} s.t. by [cor. 59.4]:

$$\mathbf{u} = \sum_{i=1}^n \alpha_i \mathbf{b}_i$$

the coefficients $\{\alpha_i\}_{i=1}^n$ need to be determined. We know that:

$$\begin{aligned}\mathbf{v} - \mathbf{u} &\perp \mathbf{b}_1, \dots, \mathbf{v} - \mathbf{u} \perp \mathbf{b}_n \\ \Rightarrow \left(\mathbf{v} - \sum_{i=1}^n \alpha_i \mathbf{b}_i \right) \cdot \mathbf{b}_j &= 0 \quad j = 1, \dots, n\end{aligned}$$

this linear system of equations can be rewritten as:

$$(\mathbf{b}_1, \dots, \mathbf{b}_n) \begin{pmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_n \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_n \end{pmatrix} \mathbf{v}$$

Proof 59.20: Corollary 59.27

Let $\mathbf{Q}\Lambda\mathbf{Q}^\top$ be the eigendecomposition^[cor. 59.12] of \mathbf{A} then it follows:

$$\begin{aligned}\min_{\|\mathbf{n}\|=1} \mathbf{n}^\top \mathbf{A} \mathbf{n} &= \min_{\|\mathbf{n}\|=1} \mathbf{n}^\top (\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{n} \\ &= \min_{\|\mathbf{n}\|=1} (\mathbf{Q}^\top \mathbf{n})^\top \Lambda (\mathbf{Q}^\top \mathbf{n}) \\ &= \min_{\|\mathbf{x}\|=1} \mathbf{x}^\top \Lambda \mathbf{x} \quad \mathbf{x} := \mathbf{Q}^\top \mathbf{n} \\ &= \min_{\|\mathbf{x}\|=1} \sum_{i=1}^n x_i^2 \Lambda_{ii} = \min_{\|\mathbf{x}\|=1} \sum_{i=1}^n x_i^2 \lambda_i\end{aligned}$$

Thus in order to obtain the minimum value we need to choose the eigenvector that leads to the smallest eigenvalue.

5. Norms

Proof 59.21: ?? 59.21

$$|\mathbf{u} \cdot \mathbf{v}| \stackrel{\text{eq. (59.81)}}{=} \|\mathbf{u}\| \|\mathbf{v}\| |\cos \theta| \leq \|\mathbf{u}\| \|\mathbf{v}\|$$

Proof 59.22: Definition 59.61

$$\begin{aligned}\|\mathbf{u} + \mathbf{v}\|^2 &= (\mathbf{u} + \mathbf{v})(\mathbf{u} + \mathbf{v}) = \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 + 2(\mathbf{u} \cdot \mathbf{v}) \\ \text{from cauchy schwartz we know:} \quad &\mathbf{u} \cdot \mathbf{v} \leq \|\mathbf{u} \cdot \mathbf{v}\| \stackrel{\text{eq. (59.91)}}{\leq} \|\mathbf{u}\| \|\mathbf{v}\| \\ \|\mathbf{u} + \mathbf{v}\|^2 &\leq \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 + 2(\|\mathbf{u}\| \|\mathbf{v}\|) = (\|\mathbf{u}\| + \|\mathbf{v}\|)^2\end{aligned}$$

6. Decompositions

19.6.1. Symmetric - Antisymmetric

Definition 59.88 Symmetric - Antisymmetric Decomposition: Any matrix $\mathbf{A} \in \mathbb{K}^{n \times n}$ can be decomposed into the sum of a *symmetric matrix*^[def. 59.23] \mathbf{A}^{sym} and a *skewsymmetric matrix*? \mathbf{A}^{skew} :

$$\begin{aligned}\mathbf{A} &= \mathbf{A}^{\text{sym}} + \mathbf{A}^{\text{skew}} \\ \mathbf{A}^{\text{sym}} &= \frac{1}{2} \left(\mathbf{A} + \mathbf{A}^\top \right) \\ \mathbf{A}^{\text{skew}} &= \frac{1}{2} \left(\mathbf{A} - \mathbf{A}^\top \right)\end{aligned}\quad (59.138)$$

19.6.2. SVD

Proof 59.23 [Corollary 59.5]: $\mathbf{B} := \mathbf{A}^\top \mathbf{A}$ corresponds to a *symmetric positive definite form*^[def. 59.75]:

$$\mathbf{x}^\top \mathbf{B} \mathbf{x} = \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} = \|\mathbf{Ax}\|_2^2 > 0$$

thus Proposition 59.6 follows immediately from [Corollary 59.2].

Proof 59.24 Proposition 59.6:

$$\begin{aligned}\mathbf{A}^\top \mathbf{A} &\stackrel{\text{SVD}}{=} \left(\mathbf{U} \Sigma \mathbf{V}^\top \right)^\top \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{V} \Sigma^\top \underbrace{\mathbf{U}^\top \mathbf{U}}_{\mathbf{I}_m} \Sigma \mathbf{V}^\top = \mathbf{V} \Sigma^\top \Sigma \mathbf{V}^\top \\ &\implies \mathbf{V}^\top \mathbf{A}^\top \mathbf{A} \mathbf{V} = \Sigma^\top \Sigma\end{aligned}$$

19.6.3. Eigendecomposition

Proof 59.25 Definition 59.86:

$$\mathbf{AX} = [\lambda_1 \mathbf{x}_1 \dots \lambda_n \mathbf{x}_n] = \mathbf{X}\Lambda$$

Geometry

Corollary 60.1 Affine Transformation in 1D: Given: numbers $x \in \Omega$ with $\Omega = [a, b]$. The affine transformation of $\phi : \Omega \rightarrow \Omega$ with $y \in \Omega = [c, d]$ is defined by:

$$y = \phi(x) = \frac{d - c}{b - a} (x - a) + c \quad (60.1)$$

Proof 60.1: [cor. 60.1] By [def. 59.45] we want a function $f : [a, b] \rightarrow [c, d]$ that satisfies:

$$f(a) = c \quad \text{and} \quad f(b) = d$$

additionally $f(x)$ has to be a linear function ([def. 54.15]), that is the output scales the same way as the input scales.

Thus it follows:

$$\frac{d - c}{b - a} = \frac{f(x) - f(a)}{x - a} \iff f(x) = \frac{d - c}{b - a} (x - a) + c$$

Trigonometry

1. Trigonometric Functions

0.1.1. Sine

Definition 60.1 Sine:

$$\sin \alpha = \frac{\text{opposite}}{\text{hypotenuse}} = \frac{a}{c} \quad (60.2)$$

0.1.2. Cosine

Definition 60.2 Cosine:

$$\cos \alpha = \frac{\text{adjacent}}{\text{hypotenuse}} = \frac{b}{c} \quad (60.3)$$

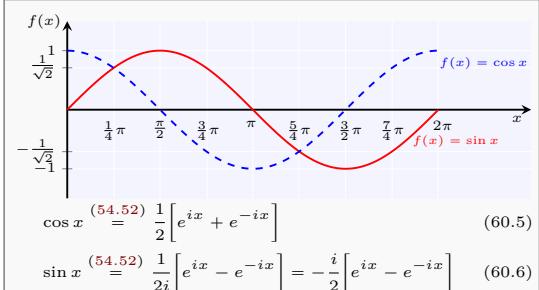
0.1.3. Tangens

Definition 60.3 Tangens:

$$\tan \alpha = \frac{\text{opposite}}{\text{adjacent}} = \frac{a}{b} = \frac{a/c}{b/c} = \frac{\sin \alpha}{\cos \alpha} \quad (60.4)$$

0.1.4. Trigonometric Functions and the Unit Circle

Sine and Cosine



Note

Using theorem 60.1 if follows:
 $\cos(\alpha \pm \pi) = -\cos \alpha$ and $\sin(\alpha \pm \pi) = -\sin \alpha$ (60.7)

0.1.5. Sinh

Definition 60.4 Sinh:

$$\sinh x \stackrel{(eq. (54.52))}{=} \frac{1}{2} [e^x - e^{-x}] = -i \sin(ix) \quad (60.8)$$

Property 60.1: $\sinh x = 0$ has a unique root at $x = 0$.

0.1.6. Cosh

Definition 60.5 Cosh:

$$\cosh x \stackrel{(54.52)}{=} \frac{1}{2} [e^x + e^{-x}] = \cos(ix) \quad (60.9)$$

$$(60.10)$$

Property 60.2: $\cosh x$ is strictly positive.

Proof 60.2:
 $e^x = \cosh x + \sinh x$ $e^{-x} = \cosh x - \sinh x$ (60.11)

2. Addition Theorems

Theorem 60.1 Addition Theorems:

$$\sin(\alpha \pm \beta) = \sin \alpha \cos \beta \pm \cos \alpha \sin \beta \quad (60.12)$$

$$\cos(\alpha \pm \beta) = \cos \alpha \cos \beta \mp \sin \alpha \sin \beta \quad (60.13)$$

3. Werner Formulas

Werner Formulas

$$\sin \alpha \cos \beta = \frac{1}{2} [\sin(\alpha + \beta) + \sin(\alpha - \beta)] \quad (60.14)$$

$$\sin \alpha \sin \beta = \frac{1}{2} [\cos(\alpha - \beta) - \cos(\alpha + \beta)] \quad (60.15)$$

$$\cos \alpha \cos \beta = \frac{1}{2} [\cos(\alpha + \beta) + \cos(\alpha - \beta)] \quad (60.16)$$

Note

Using theorem 60.1 if follows:
 $\cos(\alpha \pm \pi) = -\cos \alpha$ and $\sin(\alpha \pm \pi) = -\sin \alpha$ (60.17)

4. Law of Cosines

Law 60.1 Law of Cosines

[proof 60.3]: relates the three side of a general triangle to each other.

$$a^2 = b^2 + c^2 - 2bc \cos \theta_{b,c} \quad (60.18)$$

Law 60.2 Law of Cosines for Vectors

[proof 60.4]: relates the length of vectors to each other.

$$\|\mathbf{a}\|^2 = \|\mathbf{c} - \mathbf{b}\|^2 = \|\mathbf{b}\|^2 + \|\mathbf{c}\|^2 - 2\|\mathbf{b}\|\|\mathbf{c}\| \cos \theta_{\mathbf{b},\mathbf{c}} \quad (60.19)$$

Law 60.3 Pythagorean theorem: special case of ?? for right triangle:

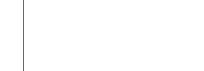
$$a^2 = b^2 + c^2 \quad (60.20)$$

1. Proofs

Proof 60.3: Law 60.1 From the definition of the sine and cosine we know that:

$$\sin \theta = \frac{h}{b} \Rightarrow h \quad \text{and} \quad \cos \theta = \frac{d}{b}$$

$$\begin{aligned} \frac{e}{a} &= c - d = c - b \cos \theta \\ \frac{e^2}{a^2} &= \frac{e^2}{c^2} + \frac{h^2}{b^2} = c^2 - 2cb \cos \theta + b^2 \cos^2 \theta + b^2 \sin^2 \theta \\ &= c^2 + b^2 - 2bc \cos \theta \end{aligned}$$



Proof 60.4: Law 60.2 Notice that $\mathbf{c} = \mathbf{a} + \mathbf{b} \Rightarrow \mathbf{a} = \mathbf{c} - \mathbf{b}$ and we can either use ?? 60.3 or notice that:

$$\begin{aligned} \|\mathbf{c} - \mathbf{b}\|^2 &= (\mathbf{c} - \mathbf{b}) \cdot (\mathbf{c} - \mathbf{b}) \\ &= \mathbf{c} \cdot \mathbf{c} - 2\mathbf{c} \cdot \mathbf{b} + \mathbf{b} \cdot \mathbf{b} \\ &= \|\mathbf{c}\|^2 + \|\mathbf{b}\|^2 - 2(\|\mathbf{c}\|\|\mathbf{b}\| \cos \theta) \end{aligned}$$

Topology

Definition 61.1 Topology of set τ : Let X be a set. A collection τ of open?? subsets of X is called topology of X if it satisfies:

- $\emptyset \in \tau$ and $X \in \tau$
- Any finite or infinite union of subsets of τ is contained in τ :

$$\{U_i : i \in I\} \subseteq \tau \implies \bigcup_{i \in I} U_i \in \tau \quad (61.1)$$

- The intersection of a finite number of elements of τ also belongs to τ :

$$\{U_i\}_{i=1}^n \in \tau \implies U_1 \cap \dots \cap U_n \in \tau \quad (61.2)$$

Definition 61.2 Topological Space (X, τ) : Is an ordered pair (X, τ) , where X is a set and τ is a topology^[def. 61.1] on X .

Numerical Methods

Machine Arithmetic's

Machine/Floating Point Numbers

Definition 62.1 (IEEE)

Institute of Electrical and Electronics Engineers:

Is a engineering associations that defines a standard on how computers should treat machine numbers in order to have certain guarantees.

Definition 62.2 Machine/Floating Point Numbers

M : Computers are only capable to represent a finite, discrete set of the real numbers $M \subset \mathbb{R}$

1.1.1. Floating Point Arithmetic's

$$x\tilde{\Omega}y = f(x\Omega y)$$

Corollary 62.1 Closure:

Machine numbers \mathbb{F} are not closed^[def. 50.7] under basic arithmetic operations:

$$\mathbb{F}\Omega\mathbb{F} \mapsto \mathbb{F} \quad \Omega = \{+, -, *, /\} \quad (62.1)$$

Note

Corollary 62.1 provides a problem as the computer can only represent floating point number \mathbb{F} .

Definition 62.3 Overflow: Result is bigger then the biggest representable floating point number.

Definition 62.4 Underflow: Result is smaller then the smaller representable floating point number i.e. to close to zero.

1.1.2. The Rounding Unit

Definition 62.5

Rounding Function/Unit

rd/\tilde{r} :

Let $x \in \mathbb{K}$ be a number real or complex number. The rounding function approximates x by the nearest machine number $\tilde{x} \in \mathbb{F}$:

$$rd : \begin{cases} \mathbb{R} \mapsto \mathbb{F} \\ x \mapsto \max_{\tilde{x} \in \mathbb{F}} \arg \min |x - \tilde{x}| \end{cases} \quad (62.2)$$

Notes

- If this is ambiguous (there are two possibilities), then it takes the larger one.
- Basic arithmetic rules such as associativity do no longer hold for operations such as addition and subtraction.

Definition 62.6 Floating Point Operation

$\tilde{\Omega}$:

Is a basic arithmetic operation between two floating point numbers $x \in \mathbb{F}$ rounded back to the nearest floating point number:

$$\mathbb{F}\tilde{\Omega}\mathbb{F} \mapsto \mathbb{F} \quad \tilde{\Omega} := rd \circ \Omega \quad (62.3) \\ \Omega = \{+, -, *, /\}$$

Definition 62.7 Absolute Error: Let $\tilde{x} \in \mathbb{K}$ be an approximation of $x \in \mathbb{K}$ then the absolute error is defined by:

$$\epsilon_{abs} := |x - \tilde{x}| \quad (62.4)$$

Definition 62.8 Relative Error: Let $\tilde{x} \in \mathbb{K}$ be an approximation of $x \in \mathbb{K}$ then the relative error is defined by:

$$\epsilon_{abs} := \frac{|x - \tilde{x}|}{|x|} \quad (62.5)$$

Note

We are interested in the relative error as it controls the number of correct/significant digits l of the approximation \tilde{x} of $x \in \mathbb{K}$:

$$\epsilon_{abs} := \frac{|x - \tilde{x}|}{|x|} \leqslant 10^l \quad l \in \mathbb{N}_{>0} \quad (62.6)$$

1.1.3. The Machine Epsilon

Definition 62.9

EPS

The Machine Epsilon:

The machine epsilon EPS is the largest possible relative rounding error^[def. 62.8]:

$$EPS := \max_{x \in I \setminus 0} \frac{|rd(x) - x|}{|x|} \quad I := [\min|\mathbb{M}|, \max|\mathbb{M}|] \in \mathbb{K} \quad (62.7)$$

Corollary 62.2 Relative Error of Flop:

The relative error^[def. 62.8] of any floating point operation^[def. 62.6] is bounded by the machine epsilon^[def. 62.9]:

$$\text{EPSrel}(\tilde{\Omega}(x, y)) := \frac{|\tilde{\Omega}(x, y) - \Omega(x, y)|}{|\Omega(x, y)|} \\ = \frac{|(rd - I)\Omega(x, y)|}{|\Omega(x, y)|} \leqslant EPS \quad (62.8)$$

Corollary 62.3 EPS for Machine Number: For machine numbers EPS can be computed by:

$$EPS = \frac{1}{2}B^{1-m} \quad (62.9)$$

Type	EPS
double	$2.2 \cdot 10^{-16}$
float	$1.1 \cdot 10^{-23}$
FP16	$9.76 \cdot 10^{-4}$

Axiom of Round off Analysis

Axiom 62.1 Axiom of Round off Analysis:

Let $x, y \in \mathbb{F}$ be (normalized) floats and assume that $x\tilde{\Omega}y \in \mathbb{F}$ (i.e. no over/underflow). Then it holds that:

$$x\tilde{\Omega}y = (x\Omega y)(1 + \delta) \quad \Omega = \{+, -, *, /\} \\ \tilde{f}(x) = f(x)(1 + \delta) \quad f \in \{\exp, \sin, \cos, \log, \dots\} \quad (62.10)$$

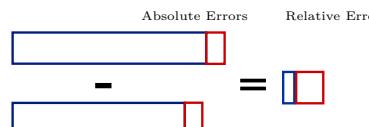
with $|\delta| < EPS$

Explanation 62.1 (axiom 62.1). gives us a guarantee that for any two floating point numbers $x, y \in \mathbb{F}$, any operation involving them will give a floating point result which is within a factor of $1 + \delta$ of the true result $x\Omega y$.

1.1.4. Cancellation

Definition 62.10 Cancellation:

Is the extreme amplification of relative errors^[def. 62.8] when subtracting numbers of almost equal size.



Roundoff Errors

2.0.1. Tricks

Log-Sum-Exp Trick

The sum exponential trick is a trick that helps to calculate the log-sum-exponential in a robust way by avoiding over/underflow. The log-sum-exponential^[def. 62.11] is an expression that arises frequently in machine learning i.e. for the cross entropy loss or for calculating the evidence of a posterior prediction.

The root of the problem is that we need to calculate the exponential $\exp(x)$, this comes with two different problems:

- If x is large (i.e. 89 for single precision floats) then $\exp(x)$ will lead to overflow
- If x is very negative $\exp(x)$ will lead to underflow/0. This is not necessarily a problem but if $\exp(x)$ occurs in the denominator or the logarithm for example this is catastrophic.

Definition 62.11 Log sum Exponential:

$$\text{LogSumExp}(x_1, \dots, x_n) := \log \left(\sum_{i=1}^n e^{x_i} \right) \quad (62.11)$$

Asymptotic Complexity

1. O-Notation

3.1.1. Small $o(\cdot)$ Notation

Definition 62.13 Little o Notation:

$$f(n) = o(g(n)) \iff \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0 \quad (62.13)$$

3.1.2. Big $\mathcal{O}(\cdot)$ Notation

2. Basic Operations

4. Rate Of Convergence

Definition 62.14 Rate of Convergence: Is a way to measure the rate of convergence of a sequence $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$ to a value to \mathbf{x}^* . Let $p \in [0, 1]$ be the *rate of convergence* and define:

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}^k - \mathbf{x}^*\|} = p \quad (62.14)$$

$$\Leftrightarrow \lim_{k \rightarrow \infty} \|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq p \|\mathbf{x}^k - \mathbf{x}^*\| \quad \forall k \in \mathbb{N}_0$$

Definition 62.15 Linear/Exponential Convergence:

A sequence $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$ converges *linearly* to \mathbf{x}^* if in the asymptotic limit $k \rightarrow \infty$ if it satisfies:

$$p \in (0, 1) \quad \forall k \in \mathbb{N}_0 \quad (62.15)$$

Definition 62.16 Superlinear Convergence:

A sequence $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$ converges *superlinear* to \mathbf{x}^* if in the asymptotic limit $k \rightarrow \infty$ if it satisfies:

$$p = 1 \quad (62.16)$$

Definition 62.17 Sublinear Convergence:

A sequence $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$ converges *sublinear* to \mathbf{x}^* if in the asymptotic limit $k \rightarrow \infty$ if it satisfies:

$$p = 0 \quad \Leftrightarrow \quad \|\mathbf{x}^{k+1} - \mathbf{x}^*\| = o\left(\|\mathbf{x}^k - \mathbf{x}^*\|\right) \quad (62.17)$$

Definition 62.18 Logarithmic Convergence:

A sequence $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$ converges *logarithmically* to \mathbf{x}^* if it converges *sublinear*^[def. 62.17] and additioinally satisfies

$$p = 0 \quad \Leftrightarrow \quad \|\mathbf{x}^{k+2} - \mathbf{x}^{k+1}\| = o\left(\|\mathbf{x}^{k+1} - \mathbf{x}^k\|\right) \quad (62.18)$$

Exponential Convergence

Linear convergence is sometimes called exponential convergence. This is due to the fact that:

- We often have expressions of the form:

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \underbrace{(1 - \alpha)}_{:= p} \|\mathbf{x}^k - \mathbf{x}^*\|$$

- and that $(1 - \alpha) = \exp(-\alpha)$ from which follows that:

$$\text{eq. (62.19)} \quad \Leftrightarrow \quad \|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq e^{-\alpha} \|\mathbf{x}^k - \mathbf{x}^*\|$$

Definition 62.19 Convergence of order p : In order to distinguish *superlinear convergence* we define the order of convergence.

A sequence $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$ converges superlinear with order $p \in \{2, \dots\}$ to \mathbf{x}^* if it satisfies:

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}^k - \mathbf{x}^*\|^p} = C \quad C < 1 \quad (62.19)$$

Definition 62.20 Exponential Convergence: A sequence $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$ converges exponentially with rate p to \mathbf{x}^* if in the asymptotic limit $k \rightarrow \infty$ it satisfies:

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq p^k \|\mathbf{x}^k - \mathbf{x}^*\| \quad p < 1 \quad (62.20)$$

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \in o\left(\|\mathbf{x}^k - \mathbf{x}^*\|\right) \quad (62.21)$$

5. Basic Operations

Operation	#mul/div	#add/sub	asympt. comp
Dot Prod.	n	$n - 1$	$\mathcal{O}(n)$
Tensor Prod.	nm	0	$\mathcal{O}(nm)$
Matrix Prod.	$m nk$	$mk(n - 1)$	$\mathcal{O}(nmk)$

Linear Systems of Equations

1. Direct Methods

6.1.1. Gaussian Elimination

Definition 62.21 Pivot Elements $a_{11}, a_{22}, \dots, a_{nn}$: Are the diagonal elements of $\mathbf{A} \in \mathbb{R}^{n,n}$ that we use to zero out the column below.

Definition 62.22 Row Echelon Matrix: Is a rectangular matrix where:

- All non-zero rows are above any zero rows.
- Each pivot of a row has a larger column index than the pivot of the row above.
- All entries below a pivot are zero.

Corollary 62.4 Reduced Form Row Echelon Matrix: Is an echelon matrix^[def. 62.22] where:

- The leading entry in each non-zero row equals 1.
- Each leading one is the only entry in its column.

Note

In case of square matrix this is a unit diagonal matrix.

Definition 62.23 197 CE
Gaussian Elimination $\mathbf{A} \in \mathbb{R}^{n,n}, \mathcal{O}(n^3)$:

Is an algorithm to solve linear systems of equations:

$$\mathbf{Ax} = \mathbf{b} \iff \begin{array}{l} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{array}$$

and consists of two steps:

(1) Forward Elimination $\mathcal{O}(n^3)$ – transforming \mathbf{A} into an upper diagonal form $[\mathbf{U}|\mathbf{b}']$:

$$\begin{array}{l} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n = b_1 \\ a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n = b_2 \\ a_{33}x_3 + \dots + a_{3n}x_n = b_3 \\ \vdots \\ a_{nn}x_n = b_n \end{array}$$

(2) Back Substitution Elimination $\mathcal{O}(n^2)$ – calculating the unknown's \mathbf{x} from \mathbf{U} :

Gauss Jordan Elimination

Is in principle the same as Gauss elimination but reduce the matrix into row-reduced echelon form^[def. 62.22].

Forward Elimination

Algorithm 62.1 Forward Elimination:

Transforms $\mathbf{Ax} = \mathbf{b}$ into row-echelon form^[def. 62.22]:

Given:

```
1: for  $k = 1, \dots, n - 1$  do
2:   pivot  $\leftarrow \mathbf{A}(k, k)$ 
3:   for  $i = k + 1, \dots, n$  do
4:      $l_{ik} \leftarrow \frac{\mathbf{A}(i, k)}{\text{pivot}}$ 
5:     for  $j = k + 1, \dots, n$  do
6:        $a_{ij} \leftarrow \mathbf{A}(i, j) = \mathbf{A}(i, j) - l_{ik} \mathbf{A}(k, j)$ 
7:     end for
8:   end for
9: end for
```

Corollary 62.5 Complexity:

$$\sum_{i=1}^{n-1} (n-1)(2(n-i)+3) = n(n-1) \left(\frac{2}{3}n + \frac{7}{6} \right) = \mathcal{O}\left(\frac{2}{3}n^3\right)$$
(62.22)

Backward Substitution

Algorithm 62.2 Backward Substitution:

Given \mathbf{U} :

```
1:  $x_n \leftarrow \frac{b_n}{a_{nn}}$ 
2: for  $i = n - 1, n - 2, \dots, 1$  do
3:    $x_i = \frac{b_i - \sum_{j=i+1}^n a_{ij} x_j}{a_{ii}}$ 
4: end for
```

Corollary 62.6 Complexity:

$$\sum_{i=1}^{n-1} 2(n-i) + 1 = \mathcal{O}(n^2)$$
(62.23)

By Rank-1 Modifications

6.1.2. LU-Decomposition

Definition 62.24

LU Decomposition

$\mathcal{O}(n^3)$:

Decomposes a matrix \mathbf{A} in an upper and lower triangular part in order to solve a system of linear equations.

Given: $\mathbf{PA} = \mathbf{LU}$ we can compute:

- ① $\mathbf{Ly} = \mathbf{Pb}$
- ② $\mathbf{Ux} = \mathbf{y}$

Corollary 62.7

[proof ??]

LU decomposition Complexity:

$$\frac{2}{3}n^3 + \frac{1}{3}n^2$$

Solving Multiple Systems of Equations

6.1.3. Symmetric Matrices

LDL-Decomposition

6.1.4. Symmetric Positive Definite Matrices

For linear systems with s.p.d.^[def. 59.75] matrices \mathbf{A} the LU-decomposition^[def. 62.24] simplifies to the Cholesky Decomposition^[def. 62.25].

Cholesky Decomposition

Definition 62.25

Cholesky Decomposition

$\frac{1}{3}\mathcal{O}(n^3)$:

Let \mathbf{A} be a s.p.d.^[def. 59.75] then it can be factorized into:

$$\mathbf{A} = \mathbf{GG}^\top \quad \text{with} \quad \mathbf{G} := \mathbf{LD}^{1/2}$$
(62.24)

Corollary 62.8

[proof 62.5]

Cholesky decomposition Complexity:

$$\frac{1}{3}n^3 + \frac{1}{3}n^2$$

2. Iterative Methods

7. Non-linear Systems of Equations

1. Iterative Methods

Definition 62.26

General Non-linear System of Equations (NLSE) F :
Is a system of non-linear equations F (that do **not** satisfy linearity??):
 $F : \subseteq \mathbb{R}^n \mapsto \mathbb{R}^n$ seek to find $\mathbf{x} \in \mathbb{R}^n : F(\mathbf{x}) = \mathbf{0}$ (62.25)

Definition 62.27 Stationary m -point Iteration ϕ_F :

Let $n, m \in \mathbb{R}$ and let $U \subseteq (\mathbb{R}^n)^m = \mathbb{R}^n \times \dots \times \mathbb{R}^n$ be a set.
A function $\phi : U \mapsto \mathbb{R}^n$, is called (m -point) iteration function if it produces an iterative sequence $(\mathbf{x}^{(k)})_k$ of approximate solutions to eq. (62.25), using the m most recent iterates:

$$\mathbf{x}^{(k)} = \phi_F(\mathbf{x}^{(k-1)}, \dots, \mathbf{x}^{(k-m)}) \quad (62.26)$$

Initial Guess $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(m-1)}$

Note

Stationary as ϕ does no explicitly depend on k .

Definition 62.28 Fixed Point \mathbf{x}^* :

Is a point \mathbf{x}^* for which the sequence does not change anymore:

$$\mathbf{x}^{(k-1)} = \mathbf{x}^*$$

$$\mathbf{x}^* = \phi_F(\mathbf{x}^{(k-1)}, \dots, \mathbf{x}^{(k-m)}) \quad \text{with} \quad \vdots$$

$$\mathbf{x}^{(k-m)} = \mathbf{x}^*$$

$$(62.27)$$

7.1.1. Convergence

Question

Does the sequence $(\mathbf{x}^{(k)})_k$ converge to a limit:

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^* \quad (62.28)$$

7.1.2. Consistency

Definition 62.29 Consistent m -point Iterative Method:

A stationary m -point method^[def. 62.27] is *consistent* with a non-linear system of equations^[def. 62.26] F iff:

$$F(\mathbf{x}^*) \iff \phi_F(\mathbf{x}^*, \dots, \mathbf{x}^*) = \mathbf{x}^* \quad (62.29)$$

7.1.3. Speed of Convergence

2. Fixed Point Iterations $m = 1$

Definition 62.30 Fixed Point Iteration: Is a 1-point method $\phi_F : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ that seeks a fixed point \mathbf{x}^* to solve $F(\mathbf{x}) = 0$:

$$\mathbf{x}^{(k+1)} = \phi_F(\mathbf{x}^{(k)}) \quad \text{Initial Guess: } \mathbf{x}^{(0)} \quad (62.30)$$

Corollary 62.9 Consistency: If ϕ_F is *continuous* and $\mathbf{x}^* = \lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$ then \mathbf{x}^* is a fixed point^[def. 62.28] of ϕ .

Algorithm 62.3 Fixed Point Iteration:

Input: Initial Guess: $\mathbf{x}^{(0)}$

1: Rewrite $F(\mathbf{x}) = 0$ into a form of $\mathbf{x} = \phi_F(\mathbf{x})$
 \triangleright There exist many ways

2: for $k = 1, \dots, T$ do

3: Use the fixed point method:

$$\mathbf{x}^{(k+1)} = \phi_F(\mathbf{x}^{(k)}) \quad (62.31)$$

4: end for

8. Numerical Quadrature

Definition 62.31 Order of a Quadrature Rule:

The order of a quadrature rule $\mathcal{Q}_n : C^0([a, b]) \rightarrow \mathbb{R}$ is defined as:

$$\text{order}(\mathcal{Q}_n) := \max \left\{ n \in \mathbb{N}_0 : \mathcal{Q}_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_n \right\} + 1 \quad (62.32)$$

Thus it is the maximal degree+1 of polynomials (of degree maximal degree) \mathcal{P} maximal degree for which the quadrature rule yields exact results.

Note

Is a quality measure for quadrature rules.

1. Composite Quadrature

Definition 62.32 Composite Quadrature:

Given a mesh $\mathcal{M} = \{a = x_0 < x_1 < \dots < x_m = b\}$ apply a Q.R. \mathcal{Q}_n to each of the mesh cells $I_j := [x_{j-1}, x_j] \quad \forall j = 1, \dots, m \triangleq \text{p.w. Quadrature:}$

$$\int_a^b f(t) dt = \sum_{j=1}^m \int_{x_{j-1}}^{x_j} f(t) dt = \sum_{j=1}^m \mathcal{Q}_n(f|_{I_j}) \quad (62.33)$$

Lemma 62.1 Error of Composite quadrature Rules:

Given a function $f \in C^k([a, b])$ with integration domain:

$$\sum_{i=1}^m h_i = |b - a| \quad \text{for } \mathcal{M} = \{x_j\}_{j=1}^m$$

Let: $h_{\mathcal{M}} = \max_j |x_j - x_{j-1}|$ be the mesh-width

Assume an equal number of quadrature nodes for each interval $I_j = [x_{j-1}, x_j]$ of the mesh \mathcal{M} i.e. $n_j = n$.

Then the error of a quadrature rule $\mathcal{Q}_n(f)$ of order q is given by:

$$\begin{aligned} \epsilon_n(f) &= \mathcal{O}(n^{-\min\{k, q\}}) = \mathcal{O}(h_{\mathcal{M}}^{\min\{k, q\}}) \quad \text{for } n \rightarrow \infty \\ [\text{cor. 54.6}] \quad \epsilon_n(f) &= \mathcal{O}(n^{-q}) = \mathcal{O}\left(h_{\mathcal{M}}^q\right) \quad \text{with } h_{\mathcal{M}} = \frac{1}{n} \end{aligned} \quad (62.34)$$

Definition 62.33 Complexity W:

Is the number of function evaluations \triangleq number of quadrature points.

$$W(\mathcal{Q}(f)_n) = \#\text{f-eval} \triangleq n \quad (62.35)$$

Lemma 62.2 Error-Complexity $W(\epsilon_n(f))$:

Relates the complexity to the quadrature error.

Assuming and quadrature error of the form :

$$\epsilon_n(f) = \mathcal{O}(n^{-q}) \iff \epsilon_n(f) = cn^{-q} \quad c \in \mathbb{R}_+$$

the error complexity is algebraic (??) and is given by:

$$W(\epsilon_n(f)) = \mathcal{O}(\epsilon_n^{1/q}) = \mathcal{O}\left(\sqrt[q]{\epsilon_n}\right) \quad (62.36)$$

Proof 62.1: lemma 62.2: Assume: we want to reduce the error by a factor of ϵ_n by increasing the number of quadrature points $n_{\text{new}} = a \cdot n_{\text{old}}$.

Question: what is the additional effort (#f-eval) needed in order to achieve this reduction in error?

$$\frac{c \cdot n_{\text{new}}^q}{c \cdot n_{\text{old}}^q} = \frac{1}{\epsilon_n} \Rightarrow n_{\text{new}} = n_{\text{old}} \cdot \sqrt[q]{\epsilon_n} = \mathcal{O}(\sqrt[q]{\epsilon_n}) \quad (62.37)$$

8.1.1. Simpson Integration

Definition 62.34 Simpson Integration:

Filtering Algorithms

10. Signals

Definition 62.35 Time Discrete Signal: Is a bounded sequence^[def. 51.2] $(x_j)_{j \in \mathbb{Z}} \in l^\infty(\mathbb{Z})$.

Definition 62.36 Sampling:

Corollary 62.10 Finite Time Discrete Signal:

11. Channels/Filters

Definition 62.37 Channel/Filter: F : Is a mapping of signals to signals $F :: l^\infty(\mathbb{Z}) \mapsto l^\infty(\mathbb{Z})$ (62.38)

Property 62.1 Finite Channel/Filter: A filter $F : l^\infty(\mathbb{Z}) \mapsto l^\infty(\mathbb{Z})$

Property 62.2 Causal Channel/Filter:

Explanation 62.3. The response cannot start before the signal has been feed into the filter.

Definition 62.38 Time Shift Operator: S_m

Property 62.3 Time-invariant Channel/Filter:

Explanation 62.4. The response of the filter should not depend at which time we pass the signal to the filter.

Property 62.4 Linear Channel/Filter:

Definition 62.39 Linear Time-invariant Finite Input Response Filter LT-FIR:

1. Impulse Responses

Definition 62.40 Impulse:

Definition 62.41 Impulse Response h :

Corollary 62.11 [proof 62.2] **Signal in terms of Impulse Responses:** We can write any arbitrary discrete signal as weighted sum of time shifted impulses:

$$F(x_j) = \quad (62.39)$$

$$(F(x_j))_j = \quad (62.40)$$

Proof 62.2 [cor. 62.11]:

2. Discrete Convolution

Definition 62.42 LT-FIR formula:

Proofs

Proof 62.3 Log Sum Trickformula 62.1:

$$\begin{aligned} \text{LSE} &= \log\left(\sum_{i=1}^n e^{x_i}\right) = \log\left(\sum_{i=1}^n e^{x_i - a} e^a\right) \\ &= \log\left(e^a \sum_{i=1}^n e^{x_i - a}\right) = \log\left(\sum_{i=1}^n e^{x_i - a}\right) + \log(e^a) \\ &= \log\left(\sum_{i=1}^n e^{x_i - a}\right) + a \end{aligned}$$

Proof 62.4 LU-Complexity^[cor. 62.7]:

For eliminating the first column we need to eliminate $n - 1$ rows by n additions and n multiplications which equals $(n - 1)2n$. For the second column we need for $n - 2$ rows $n - 1$ additions and $n - 1$ multiplications which equals $(n - 2)2(n - 1)$ thus to eliminate all n columns we have:

$$\sum_{i=1}^n (n - i + 1) \cdot 2(n - i)$$

using the index $l = n - i + 1$ we can write this as:

$$\begin{aligned} \sum_{i=1}^n (n - i + 1) \cdot 2(n - i) &= 2 \sum_{l=0}^n (j + 1) \cdot (j) = 2 \sum_{l=0}^n j^2 + 1 \\ &= 2 \left(\frac{1}{3} n^3 - \frac{1}{3} n \right) \end{aligned}$$

Proof 62.5 Cholesky Complexity^[cor. 62.8]: \mathbf{U} and \mathbf{L} "are the same" as we have a s.p.d. matrix s.t. we can simply half the forward elimination complexity of the LU-decomposition^[cor. 62.7]:

$$\frac{1}{2} \frac{2}{3} n^3 + \frac{1}{3} n^2 \quad (62.41)$$

Optimization

Definition 63.1 Fist Order Method: A first-order method is an algorithm that chooses the k -th iterate in $\mathbf{x}_0 + \text{span}\{\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_{k-1})\} \quad \forall k = 1, 2, \dots \quad (63.1)$

Note

Gradient descent is a first order method

1. Linear Optimization

1. Polyhedra

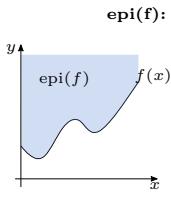
Definition 63.2 Polyhedron: Is a set $P \subseteq \mathbb{R}^n$ that can be described by the finite intersection of m closed half spaces??:

$$P = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\} = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{a}_j \mathbf{x} \leq b_j, j = 1, \dots, m\}$$

$$\mathbf{A} \in \mathbb{R}^{m \times n} \quad \mathbf{b} \in \mathbb{R}^m \quad (63.2)$$

1.1.1. Polyhedral Function

Definition 63.3 Epigraph/Subgraph



The epigraph of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is defined as the set of points that lie above its graph:

$$\text{epi}(f) := \{(\mathbf{x}, y) \in \mathbb{R}^n \mid y \geq f(\mathbf{x})\} \subseteq \mathbb{R}^{n+1} \quad (63.3)$$

Definition 63.4 Polyhedral Function: A function f is *polyhedral* if its epigraph $\text{epi}(f)^{[\text{def. 63.3}]}$ is a polyhedral set^[def. 63.2]:

$$f \text{ is polyhedral} \iff \text{epi}(f) \text{ is polyhedral} \quad (63.4)$$

2. Lagrangian Optimization Theory

Definition 63.5 (Primal) Constraint Optimization: Given an optimization problem with domain $\Omega \subseteq \mathbb{R}^d$:

$$\begin{aligned} & \min_{\mathbf{w} \in \Omega} f(\mathbf{w}) \\ \text{s.t.} \quad & g_i(\mathbf{w}) \leq 0 \quad 1 \leq i \leq k \\ & h_j(\mathbf{w}) = 0 \quad 1 \leq j \leq m \end{aligned}$$

Definition 63.6 Lagrange Function:

$$\mathcal{L}(\alpha, \beta, \mathbf{w}) := f(\mathbf{w}) + \alpha g(\mathbf{w}) + \beta h(\mathbf{w}) \quad (63.5)$$

Extremal Conditions

$$\begin{aligned} \nabla \mathcal{L}(\mathbf{x}) &= 0 && \text{Extremal point } \mathbf{x}^* \\ \frac{\partial}{\partial \beta} \mathcal{L}(\mathbf{x}) &= h(\mathbf{x}) = 0 && \text{Constraint satisfaction} \end{aligned}$$

For the inequality constraints $g(\mathbf{x}) \leq 0$ we distinguish two situations:

Case I : $g(\mathbf{x}^*) < 0$ switch const. off

Case II : $g(\mathbf{x}^*) \geq 0$ optimze using active eq. constr.

$$\frac{\partial}{\partial \alpha} \mathcal{L}(\mathbf{x}) = g(\mathbf{x}) = 0 \quad \text{Constraint satisfaction}$$

Definition 63.7 Lagrangian Dual Problem:

Is given by:

$$\text{Find } \max \theta(\alpha, \beta) = \inf_{\mathbf{w} \in \Omega} \mathcal{L}(\mathbf{w}, \alpha, \beta)$$

$$\text{s.t. } \alpha_i \geq 0 \quad 1 \leq i \leq k$$

Solution Strategy

- Find the extremal point \mathbf{w}^* of $\mathcal{L}(\mathbf{w}, \alpha, \beta)$:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^*} = 0 \quad (63.6)$$
- Insert \mathbf{w}^* into \mathcal{L} and find the extremal point β^* of the resulting dual Lagrangian $\theta(\alpha, \beta)$ for the active constraints:

$$\frac{\partial \theta}{\partial \beta} \Big|_{\beta=\beta^*} = 0 \quad (63.7)$$
- Calculate the solution $\mathbf{w}^*(\beta^*)$ of the constraint minimization problem.

Value of the Problem

Value of the problem: the value $\theta(\alpha^*, \beta^*)$ is called the value of problem (α^*, β^*) .

Theorem 63.1 Upper Bound Dual Cost: Let $\mathbf{w} \in \Omega$ be a feasible solution of the primal problem^[def. 63.5] and (α, β) a feasible solution of the respective dual problem^[def. 63.7]. Then it holds that:

$$f(\mathbf{w}) \geq \theta(\alpha, \beta) \quad (63.8)$$

Proof 63.1:

$$\begin{aligned} \theta(\alpha, \beta) &= \inf_{\mathbf{u} \in \Omega} \mathcal{L}(\mathbf{u}, \alpha, \beta) \leq \mathcal{L}(\mathbf{w}, \alpha, \beta) \\ &= f(\mathbf{w}) + \sum_{i=1}^k \underbrace{\alpha_i}_{\geq 0} g_i(\mathbf{w}) + \sum_{j=1}^m \underbrace{\beta_j}_{\leq 0} h_j(\mathbf{w}) \\ &\leq f(\mathbf{w}) \end{aligned}$$

Corollary 63.1 Duality Gap Corollary: The value of the dual problem is upper bounded by the value of the primal problem:

$$\sup \{\theta(\alpha, \beta) : \alpha \geq 0\} \leq \inf \{f(\mathbf{w}) : \mathbf{g}(\mathbf{w}) \leq 0, \mathbf{h}(\mathbf{w}) = 0\} \quad (63.9)$$

Theorem 63.2 Optimality: The triple $(\mathbf{w}^*, \alpha^*, \beta^*)$ is a saddle point of the Lagrangian function for the primal problem, if and only if its components are optimal solutions of the primal and dual problems and if there is no duality gap, that is, the primal and dual problems having the same value:

$$f(\mathbf{w}^*) = \theta(\alpha^*, \beta^*) \quad (63.10)$$

Definition 63.8 Convex Optimization: Given a convex function f and a convex set S solve:

$$\begin{aligned} & \min f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in S \end{aligned} \quad (63.11)$$

Often S is specified using linear inequalities:

$$\text{e.g. } S = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$$

Theorem 63.3 Strong Duality: Given an convex optimization problem:

$$\begin{aligned} & \min_{\mathbf{w} \in \Omega} f(\mathbf{w}) \\ \text{s.t.} \quad & g_i(\mathbf{w}) \leq 0 \quad 1 \leq i \leq k \\ & h_j(\mathbf{w}) = 0 \quad 1 \leq j \leq m \end{aligned}$$

where g_i, h_i can be written as affine functions: $y(\mathbf{w}) = \mathbf{A}\mathbf{w} - \mathbf{b}$.

Then it holds that the **duality gap** is zero and we obtain an optimal solution.

Theorem 63.4 Kuhn-Tucker Conditions: Given an optimization problem with convex domain $\Omega \subseteq \mathbb{R}^d$,

$$\begin{aligned} & \min_{\mathbf{w} \in \Omega} f(\mathbf{w}) \\ \text{s.t.} \quad & g_i(\mathbf{w}) \leq 0 \quad 1 \leq i \leq k \\ & h_j(\mathbf{w}) = 0 \quad 1 \leq j \leq m \end{aligned}$$

with $f \in C^1$ convex and g_i, h_i affine.

Necessary and sufficient conditions for a normal point \mathbf{w}^* to be an optimum are the existence of α^*, β^* s.t.:

$$\frac{\partial \mathcal{L}(\mathbf{w}, \alpha, \beta)}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^*} = 0 \quad \frac{\partial \mathcal{L}(\mathbf{w}^*, \alpha, \beta)}{\partial \beta} \Big|_{\beta=\beta^*} = 0 \quad (63.12)$$

under the conditions that:

- $\forall i_1, \dots, k \quad \alpha_i^* g_i(\mathbf{w}^*) = 0$, s.t.:
 - Inactive Constraint: $g_i(\mathbf{w}^*) < 0 \rightarrow \alpha_i = 0$.
 - Active Constraint: $g_i(\mathbf{w}^*) = 0 \rightarrow \alpha_i \geq 0$ s.t. $\alpha_i^* g_i(\mathbf{w}^*) = 0$

Consequence

We may become very sparse problems, if a lot of constraints are not active $\iff \alpha_i = 0$.

Only a few points, for which $\alpha_i > 0$ may affect the decision surface.

Combinatorics

Permutations

Definition 64.1 Permutation: A n -Permutation is the (re)arrangement of n elements of a set^[def. 50.1] \mathcal{S} of size $n = |\mathcal{S}|$ into a sequences^[def. 51.2] – **order does matter**.

Definition 64.2 Number of Permutations of a Set $n!$: Let \mathcal{S} be a set^[def. 50.1] $n = |\mathcal{S}|$ distinct objects. The number of permutations of \mathcal{S} is given by:

$$P_n(\mathcal{S}) = n! = \prod_{i=0}^{n-1} (n-i) = n \cdot (n-1) \cdot (n-2) \cdots \cdot 1 \quad (64.1)$$

Explanation 64.1. If we have i.e. three distinct elements $\{\bullet, \circ, \bullet\}$ For the first element \bullet that we arrange we have three possible choices where to put it. However this reduces the number of possible choices for the second element \bullet to only two. Consequently for the last element \bullet we have no choice left.



Definition 64.3

Number of Permutations of a Multiset:

Let \mathcal{S} be a multi set^[def. 50.3] with $n = |\mathcal{S}|$ total and k distinct objects. Let n_j be the multiplicity^[def. 50.4] of the member $j \in \{1, \dots, k\}$ of the multiset \mathcal{S} . The permutation of \mathcal{S} is given by:

$$P_{n_1, \dots, n_k}(\mathcal{S}) = \frac{n!}{n_1! \cdots n_k!} \quad \text{s.t.} \quad \sum_{j=1}^k n_j \leq n \quad k < n \quad (64.2)$$

Note

We need to divide by the permutations as sequence/order does not change if we exchange objects of the same kind (e.g. red ball by red ball) \Rightarrow less possibilities to arrange the elements uniquely.

Picking things from a bag

1. Combinations

Definition 64.4 k -Combination:

A k -combination of a set \mathcal{S} of *distinct* elements of size $n = |\mathcal{S}|$ is a subset \mathcal{S}_k (**order does not matter**) of $k = |\mathcal{S}_k|$, chosen from \mathcal{S} .

Note

Thus unlike in a permutation we just care about what we pick and not how it ends up being arranged.

Definition 64.5 Number of k -Combinations $C_{n,k}$: The number of k -combinations of a set \mathcal{S} of size $n = |\mathcal{S}|$ is given by n choose k :

$$C_{n,k} = \binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (64.3)$$

2. Variation

Definition 64.6 Variation:

A k -variation of a set \mathcal{S} of size $n = |\mathcal{S}|$ is

1. a selection/combination^[def. 64.4] of a subset \mathcal{S}_k (**order does not matter**) of k -*distinct* elements $k = |\mathcal{S}_k|$, chosen from \mathcal{S}
2. and an k arrangement/permuation^[def. 64.2] of that subset \mathcal{S}_k (with or without repetition) into a sequence^[def. 51.2]

Definition 64.7

Number of Variations without repetitions

V_k^n :

Let \mathcal{S} be a set^[def. 50.1] $n = |\mathcal{S}|$ distinct objects from which we choose k elements. The number of variations of size $k = |\mathcal{S}_k|$ of the set \mathcal{S} *without repetitions* is given by:

$$V_k^n(\mathcal{S}) = \binom{n}{k} k! = \frac{n!}{(n-k)!} \quad (64.4)$$

Note

Sometimes also denotes as P_k^n .

Definition 64.8

Number of Variations with repetitions

\bar{V}_k^n :

Let \mathcal{S} be a set^[def. 50.1] $n = |\mathcal{S}|$ distinct objects from which we choose k elements. The number of variations of size $k = |\mathcal{S}_k|$ of the set \mathcal{S} from which we *choose and always return* is given by:

$$\bar{V}_k^n(\mathcal{S}) = n^k \quad (64.5)$$

Definition 64.9 Stochastics: Is a collective term for the areas of probability theory and statistics.

Definition 64.10 Statistics: Is concerned with the analysis of data/experiments in order to draw conclusion of the underlying governing models that describe these experiments.

Definition 64.11 Probability: Is concerned with the quantification of the uncertainty of random experiments by use of statistical models. Hence it is the opposite of statistics.

Definition 64.12 Probability: Probability is the measure of the likelihood that an event will occur in a Random Experiment. Probability is quantified as a number between 0 and 1, where, loosely speaking, 0 indicates impossibility and 1 indicates certainty.

Note: Stochastics vs. Stochastic

Stochasticss is a noun and is a collective term for the areas of probability theory and statistics, while stochastic is a *adjective*, describing that a certain phenomena is governed by uncertainty i.e. a process.

Probability Theory

Definition 65.1 Probability Space $W = \{\Omega, \mathcal{F}, \mathbb{P}\}$: Is the unique triple $\{\Omega, \mathcal{F}, \mathbb{P}\}$, where Ω is its sample space, \mathcal{F} its σ -algebra of events, and \mathbb{P} its probability measure.

Definition 65.2 [example 65.1]

Sample Space Ω :

Is the set of all possible outcomes (elementary events) [cor. 65.5]) of an experiment.

Definition 65.3 [example 65.2]

Event A :

An “event” is a subset of the sample space Ω and is a property which can be observed to hold or not to hold *after* the experiment is done.

Mathematically speaking not every subset of Ω is an event and has an associated probability.

Only those subsets of Ω that are part of the corresponding σ -algebra \mathcal{F} are events and have their assigned probability.

Corollary 65.1 : If the outcome ω of an experiment is in the subset A , then the event A is said to “have occurred”.

Corollary 65.2 Complement Set A^C : Is the contrary event of A .

Corollary 65.3 The Union Set $A \cup B$: Let A, B be two events. The event “ A or B ” is interpreted as the union of both.

Corollary 65.4 The Intersection Set $A \cap B$: Let A, B be two events. The event “ A and B ” is interpreted as the intersection of both.

Corollary 65.5 The Elementary Event ω : Is a “singleton”, i.e. a subset $\{\omega\}$ containing a single outcome ω of Ω .

Corollary 65.6 The Sure Event Ω : Is equal to the sample space as it contains all possible elementary events.

Corollary 65.7 The Impossible Event \emptyset : The impossible event i.e. nothing is happening is denoted by the empty set.

Definition 65.4 The Family of All Events $\mathcal{A}/2^\Omega$: The set of all subset of the sample space Ω called family of all events is given by the power set of the sample space $\mathcal{A} = 2^\Omega$ (for finite sample spaces).

Definition 65.5 Probability

$\mathbb{P}(A)$: Is a number associated with every A , that measures the likelihood of the event to be realized “a priori”. The bigger the number the more likely the event will happen.

1. $0 \leq \mathbb{P}(A) \leq 1$
2. $\mathbb{P}(\Omega) = 1$
3. If $A \cap B = \emptyset$ then $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$

Note

We can think of the probability of an event A as the limit of the “frequency” of repeated experiments:

$$\mathbb{P}(A) = \lim_{n \rightarrow \infty} \frac{\delta_n(A)}{n} \quad \text{where} \quad \delta(A) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}$$

1. Sigma Algebras

Definition 65.6

[Proof 65.3]

Sigma Algebra

σ :

A set \mathcal{F} of subsets of Ω is called a σ -algebra on Ω if the following properties apply

- $\Omega \in \mathcal{F}$ and $\emptyset \in \mathcal{F}$
- If $A \in \mathcal{F}$ then $\Omega \setminus A = A^C \in \mathcal{F}$: The complementary subset of A is also in Ω .
- For all $A_i \in \mathcal{F}$: $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$

Explanation 65.1 ([def. 65.6]). The σ -algebra determines what events we can measure, it represents all of the possible events of the experiment that we can detect.

Thus the sigma algebra is a mathematical construct that tells us how much information we obtain once we conduct some experiment.

Corollary 65.8 $\mathcal{F}_{\min} = \mathcal{F} = \{\emptyset, \Omega\}$ is the simplest σ -algebra, telling us only if an event happened $\omega \in \Omega$ happened or not but not which one.

Corollary 65.9 $\mathcal{F}_{\max} = \mathcal{F} = 2^\Omega$ consists of all subsets of Ω and thus corresponds to full information i.e. we know if and which event happened.

Definition 65.7 Measurable Space $\{\Omega, \mathcal{F}\}$: Is the pair of a set and sigma algebra i.e. a sample space and sigma algebra $\{\Omega, \mathcal{F}\}$.

Corollary 65.10 \mathcal{F} -measurable Event $A_i \in \mathcal{F}$: The measurable events A_i of \mathcal{F} are called \mathcal{F} -measurable or measurable sets.

Definition 65.8 [Example 65.4]

Sigma Algebra generated by a subset of Ω $\sigma(\mathcal{C})$: Let \mathcal{C} be a class of subsets of Ω . The σ -algebra generated by \mathcal{C} , denoted by $\sigma(\mathcal{C})$, is the smallest sigma algebra \mathcal{F} that included all elements of \mathcal{C} .

Definition 65.9 [Example 65.5]

Borel σ -algebra $\mathcal{B}(\mathbb{R})$: The Borel σ -algebra $\mathcal{B}(\mathbb{R})$ is the smallest σ -algebra containing all open intervals in \mathbb{R} . The sets in contained in $\mathcal{B}(\mathbb{R})$ are called Borel sets.

The extension to the multi-dimensional case, $\mathcal{B}(\mathbb{R}^n)$, is straightforward.

For all real numbers $a, b \in \mathbb{R}$, $\mathcal{B}(\mathbb{R})$ contains various sets.

Why do we need Borel Sets

So far we only looked at atomic events ω , with the help of sigma algebras we are now able to measure continuous events s.a. $[0, 1]$.

Definition 65.10 Borel Set:

Corollary 65.11 Generating Borel σ -Algebra [Proof 65.1]: The Borel σ -algebra of \mathbb{R} is generated by intervals of the form $(-\infty, a]$, where $a \in \mathbb{Q}$ (\mathbb{Q} = rationals).

Definition 65.11 (\mathbb{P})-trivial Sigma Algebra: It is a σ -algebra \mathcal{F} for which each event has a probability of zero or one:

$$\mathbb{P}(A) \in \{0, 1\} \quad \forall A \in \mathcal{F} \quad (65.1)$$

Interpretation

A trivial sigma algebra means that all events are almost surely constant and that there exist no non-trivial information. An example of a trivial sigma algebra is $\mathcal{F}_{\min} = \{\Omega, \emptyset\}$.

2. Measures

Definition 65.12 Measure

μ : A measure defined on a measurable space $\{\Omega, \mathcal{F}\}$ is a function/map:

$$\mu : \mathcal{F} \mapsto [0, \infty] \quad (65.2)$$

for which holds:

- $\mu(\emptyset) = 0$
- countable additivity [def. 65.13]

Definition 65.13 Countable/ σ -Additive Function:

Given a function μ defined on a σ -algebra \mathcal{F} .

The function μ is said to be countable additive if for every countable sequence of pairwise disjoint elements $(F_i)_{i \geq 1}$ of \mathcal{F} it holds that:

$$\mu\left(\bigcup_{i=1}^{\infty} F_i\right) = \sum_{i=1}^{\infty} \mu(F_i) \quad \text{for all } F_j \cap F_k = \emptyset \quad \forall j \neq k \quad (65.3)$$

Corollary 65.12 Additive Function: A function that satisfies countable additivity, is also additive, meaning that for every $F, G \in \mathcal{F}$ it holds:

$$F \cap G = \emptyset \implies \mu(F \cup G) = \mu(F) + \mu(G) \quad (65.4)$$

Explanation 65.2. If we take two events that cannot occur simultaneously, then the probability that at least one of the events occurs is just the sum of the measures (probabilities) of the original events.

Definition 65.14 [Example 65.6]

Equivalent Measures

$\mu \sim \nu$: Let μ and ν be two measures defined on a measurable space [def. 65.7] (Ω, \mathcal{F}) . The two measures are said to be equivalent if it holds that:

$$\mu(A) > 0 \iff \nu(A) > 0 \quad \forall A \subseteq \Omega \quad (65.5)$$

this is equivalent to μ and ν having equivalent null sets:

$$\begin{aligned} \mathcal{N}_\mu &= \{A \in \mathcal{A} | \mu(A) = 0\} \\ \mathcal{N}_\nu &= \{A \in \mathcal{A} | \nu(A) = 0\} \end{aligned} \quad (65.6)$$

Definition 65.15 Measure Space $\{\mathcal{F}, \Omega, \mu\}$: The triplet of sample space, sigma algebra and a measure is called a measure space.

1. Borel Measures

Definition 65.16 Borel Measure: A Borel Measure is any measure [def. 65.12] μ defined on the Borel σ -algebra [def. 65.9] $\mathcal{B}(\mathbb{R})$.

2.1.1. The Lebesgue Measure

Definition 65.17 Lebesgue Measure on \mathbb{R} λ : Is the Borel measure [def. 65.16] defined on the measurable space $\{\mathbb{R}, \mathcal{B}(\mathbb{R})\}$ which assigns for every half-open interval $(a, b]$ its length:

$$\lambda((a, b]) := b - a \quad (65.7)$$

Corollary 65.13 Lebesgue Measure of Atomites:

- The Lebesgue measure of a set containing only one point must be zero:

$$\lambda(\{a\}) = 0 \quad (65.8)$$

- The Lebesgue measure of a set containing countably many points $A = \{a_1, a_2, \dots, a_n\}$ must be zero:

$$\lambda(A) + \sum_{i=1}^n \lambda(\{a_i\}) = 0 \quad (65.9)$$

- The Lebesgue measure of a set containing uncountably many points $A = \{a_1, a_2, \dots\}$ can be either zero, positive and finite or infinite.

3. Probability/Kolmogorov's Axioms

1931

One problem we are still having is the range of μ , by standardizing the measure we obtain a well defined measure of events.

Axiom 65.1 Non-negativity: The probability of an event is a non-negative real number:
If $A \in \mathcal{F}$ then $\mathbb{P}(A) \geq 0$ (65.10)

Axiom 65.2 Unitarity: The probability that at least one of the elementary events in the entire sample space Ω will occur is equal to one:
The certain event $\mathbb{P}(\Omega) = 1$ (65.11)

Axiom 65.3 σ -additivity: If $A_1, A_2, A_3, \dots \in \mathcal{F}$ are mutually disjoint, then:

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mathbb{P}(A_i) \quad (65.12)$$

Corollary 65.14 : As a consequence of this it follows:
 $\mathbb{P}(\emptyset) = 0$ (65.13)

Corollary 65.15 Complementary Probability:
 $\mathbb{P}(A^C) = 1 - \mathbb{P}(A)$ with $A^C = \Omega - A$ (65.14)

Definition 65.18 Probability Measure \mathbb{P} : a probability measure is function $\mathbb{P} : \mathcal{F} \mapsto [0, 1]$ defined on a σ -algebra \mathcal{F} of a sample space Ω that satisfies the probability axioms.

4. Conditional Probability

Definition 65.19 Conditional Probability: Let A, B be events, with $\mathbb{P}(B) \neq 0$. Then the conditional probability of the event A given B is defined as:

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \quad \mathbb{P}(B) \neq 0 \quad (65.15)$$

5. Independent Events

Theorem 65.1

Independent Events: Let A, B be two events. A and B are said to be independent iff:

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B) \quad \mathbb{P}(B|A) = \mathbb{P}(B), \quad \mathbb{P}(A) > 0 \quad (65.16)$$

Note

The requirement of no impossible events follows from [def. 65.19]

Corollary 65.16 Pairwise Independent Evenest: A finite set of events $\{A_i\}_{i=1}^n \in \mathcal{A}$ is pairwise independent if every pair of events is independent:

$$\mathbb{P}(A_i \cap A_j) = \mathbb{P}(A_i) \cap \mathbb{P}(A_j) \quad i \neq j, \quad \forall i, j \in \mathcal{A} \quad (65.17)$$

Corollary 65.17 Mutual Independent Evenest: A finite set of events $\{A_i\}_{i=1}^n \in \mathcal{A}$ is mutual independent if every event A_j is independent of any intersection of the other events:

$$\mathbb{P}\left(\bigcap_{i=1}^k B_i\right) = \prod_{i=1}^k \mathbb{P}(B_i) \quad \forall \{B_i\}_{i=1}^k \subseteq \{A_i\}_{i=1}^n \quad k \leq n, \quad \{A_i\}_{i=1}^n \in \mathcal{A} \quad (65.18)$$

6. Product Rule

Law 65.1 Product Rule: Let A, B be two events then the probability of both events occurring simultaneously is given by:

$$\mathbb{P}(A \cap B) = \mathbb{P}(B|A)\mathbb{P}(A) = \mathbb{P}(A|B)\mathbb{P}(B) \quad (65.19)$$

Law 65.2

Generalized Product Rule/Chain Rule: is the generalization of the product rule?? to n events $\{A_i\}_{i=1}^n$

$$\begin{aligned} \mathbb{P}\left(\bigcap_{i=i}^k E_i\right) &= \prod_{k=1}^n \mathbb{P}\left(E_k \mid \bigcap_{i=1}^{k-1} E_i\right) = \\ &= \mathbb{P}(E_n | E_{n-1} \cap \dots \cap E_1) \cdot \mathbb{P}(E_{n-1} | E_{n-2} \cap \dots \cap E_1) \dots \\ &\quad \dots \cdot \mathbb{P}(E_3 | E_2 \cap E_1) \mathbb{P}(E_2 | E_1) \mathbb{P}(E_1) \end{aligned} \quad (65.20)$$

7. Law of Total Probability

Definition 65.20 Complete Event Field: A complete event field $\{\mathcal{A}_i : i \in I \subseteq \mathbb{N}\}$ is a countable or finite partition of Ω that is the partitions $\{\mathcal{A}_i : i \in I \subseteq \mathbb{N}\}$ are a disjoint union of the sample space:

$$\bigcup_{i \in I} \mathcal{A}_i = \Omega \quad \mathcal{A}_i \cap \mathcal{A}_j = \emptyset \quad i \neq j, \forall i, j \in I \quad (65.21)$$

Theorem 65.2

Law of Total Probability/Partition Equation:

Let $\{\mathcal{A}_i : i \in I\}$ be a complete event field^[def. 65.20] then it holds for $B \in \mathcal{B}$:

$$\mathbb{P}(B) = \sum_{i \in I} \mathbb{P}(B | \mathcal{A}_i) \mathbb{P}(\mathcal{A}_i) \quad (65.22)$$

8. Bayes Theorem

Law 65.3 Bayes Rule: Let A, B be two events s.t. $\mathbb{P}(B) > 0$ then it holds:

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(B | A) \mathbb{P}(A)}{\mathbb{P}(B)} \quad \mathbb{P}(B) > 0 \quad (65.23)$$

follows directly from eq. (65.19).

Theorem 65.3 Bayes Theorem: Let $\{\mathcal{A}_i : i \in I\}$ be a complete event field^[def. 65.20] and $B \in \mathcal{B}$ a random event s.t. $\mathbb{P}(B) > 0$, then it holds:

$$\mathbb{P}(A_j | B) = \frac{\mathbb{P}(B | A_j) \mathbb{P}(A_j)}{\sum_{i \in I} \mathbb{P}(B | A_i) \mathbb{P}(A_i)} \quad (65.24)$$

proof ?? 65.2

Distributions on \mathbb{R}

1. Distribution Function

Definition 65.21 Distribution Function of \mathbb{P} F : The distribution function F induced by a probability measure \mathbb{P} on $(\mathbb{R}, \mathcal{B})$ is the function:

$$F(x) = \mathbb{P}((-\infty, x]) \quad (65.25)$$

Theorem 65.4 : A function F is the distribution function of a (unique) probability on $(\mathbb{R}, \mathcal{B})$ iff:

- F is non-decreasing
- F is right continuous
- $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow +\infty} F(x) = 1$

Corollary 65.18 : A probability \mathbb{P} is uniquely determined by a distribution function F .

That is if there exist another probability \mathbb{Q} s.t.

$$G(x) = \mathbb{Q}((-\infty, x])$$

and if $F = G$ then it follows $\mathbb{P} = \mathbb{Q}$.

2. Random Variables

A random variable X is a function/map that determines a quantity of interest based on the outcome $\omega \in \Omega$ of a random experiment. Thus X is not really a variable in the classical sense but a variable with respect to the outcome of an experiment. Its value is determined in two steps:

- ① The outcome of an experiment is a random quantity $\omega \in \Omega$
- ② The outcome ω determines (possibly various) quantities of interests \Leftrightarrow random variables

Thus a random variable X , defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ is a mapping from Ω into another space \mathcal{E} , usually $E = \mathbb{R}$ or $\mathcal{E} = \mathbb{R}^n$:

$$X : \Omega \mapsto \mathcal{E} \quad \omega \mapsto X(\omega)$$

Let now $E \in \mathcal{E}$ be a quantity of interest, in order to quantify its probability we need to map it back to the original sample space Ω :

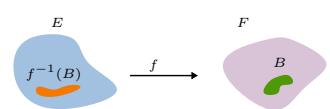
Probability for an event in Ω

$$\mathbb{P}_X(E) = \mathbb{P}(\{\omega : X(\omega) \in E\}) = \mathbb{P}(X \in E) = \overbrace{\mathbb{P}(X^{-1}(E))}^{\text{Probability for an event in } \mathcal{E}}$$

Probability for an event in E

Definition 65.22 \mathcal{E} -measurable function: Let (E, \mathcal{E}) and (F, \mathcal{F}) be two measurable spaces. A function $f : E \mapsto F$ is called measurable (relative to \mathcal{E} and \mathcal{F}) if

$$\forall B \in \mathcal{F} : f^{-1}(B) = \{\omega \in E : f(\omega) \in B\} \in \mathcal{E} \quad (65.26)$$



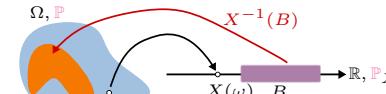
Interpretation

The pre-image^[def. 54.11] of B under f i.e. $f^{-1}(B)$ maps all values of the target space F back to the sample space Ω (for all possible $B \in \mathcal{F}$).

Definition 65.23 Random Variable: A real-valued random variable (vector) X , defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ is an \mathcal{E} -measurable function mapping, if it maps its sample space Ω into a target space (F, \mathcal{F}) :

$$X : \Omega \mapsto F \quad (\mathcal{F}^n) \quad (65.27)$$

Since X is \mathcal{E} -measurable it holds that $X^{-1} : \mathcal{F} \mapsto \mathcal{E}$



Corollary 65.19 : Usually $F = \mathbb{R}$, which usually amounts to using the Borel σ -algebra \mathcal{B} of \mathbb{R} .

Corollary 65.20 Random Variables of Borel Sets: Given that we work with Borel σ -algebras then the definition of a random variable is equivalent to (due to [cor. 65.11]):

$$\begin{aligned} X^{-1}(B) &= X^{-1}((-\infty, a]) \\ &= \{\omega \in \Omega : X(\omega) \leq a\} \in \mathcal{E} \quad \forall a \in \mathbb{R} \end{aligned} \quad (65.28)$$

Definition 65.24

Realization of a Random Variable $x = X(\omega)$: Is the value of a random variable that is actually observed after an experiment has been conducted. In order to avoid confusion lower case letters are used to indicate actual observations/realization of a random variable.

Corollary 65.21 Indicator Functions

$I_A(\omega)$: An important class of measurable functions that can be used as r.v. are indicator functions:

$$I_A(\omega) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases} \quad (65.29)$$

We know that a probability measure \mathbb{P} on \mathbb{R} is characterized by the quantities $\mathbb{P}((-\infty, a])$. Thus the quantities.

Corollary 65.22 : Let $(F, \mathcal{F}) = (\mathbb{R}, \mathcal{B})$ and let (E, \mathcal{E}) be an arbitrary measurable space. Let X be a real value function on E .

Then it holds that X is measurable if and only if

$$\begin{aligned} \{X \leq a\} &= \{\omega : X(\omega) \leq a\} = X^{-1}((-\infty, a]) \in \mathcal{E}, \forall a \in \mathbb{R} \\ \text{or} \quad \{X < a\} &\in \mathcal{E}. \end{aligned}$$

Explanation 65.3 ([cor. 65.22]). A random variable is a function that is measurable if and only if its distribution function is defined.

3. The Law of Random Variables

Definition 65.25 Law/Distribution of X

Let X be a r.v. on $(\Omega, \mathcal{F}, \mathbb{P})$, with values in (E, \mathcal{E}) , then the distribution/law of X is defined as:

$$\begin{aligned} \mathbb{P} : \mathcal{B} &\mapsto [0, 1] \quad (65.30) \\ \mathbb{P}^X(B) &= \mathbb{P}\{X \in B\} = \mathbb{P}(\omega : X(\omega) \in B) \quad \forall b \in \mathcal{B} \end{aligned}$$

Note

- Sometimes \mathbb{P}^X is also called the *image* of \mathbb{P} by X
- The law can also be written as:

$$\mathbb{P}^X(B) = \mathbb{P}(X^{-1}(B)) = (\mathbb{P} \circ X^{-1})(B)$$

Theorem 65.5 : The law/distribution of X is a probability measure \mathbb{P} on (E, \mathcal{E}) .

Definition 65.26

(Cumulative) Distribution Function

F_X : Given a real-valued r.v. then its cumulative distribution function is defined as:

$$F_X(x) = \mathbb{P}^X((-\infty, x]) = \mathbb{P}(X \leq x) \quad (65.31)$$

Corollary 65.23 : The distribution of \mathbb{P}^X of a real valued r.v. is entirely characterized by its cumulative distribution function F_X [def. 65.33].

Property 65.1:

$$\mathbb{P}(X > x) = 1 - F_X(x) \quad (65.32)$$

Property 65.2:

$$\mathbb{P}(a < X \leq b) = F_X(b) - F_X(a) \quad (65.33)$$

4. Probability Density Function

Definition 65.27 Continuous Random Variable: Is a r.v. for which a probability density function f_X exists.

Definition 65.28 Probability Density Function: Let X be a r.v. with associated cdf F_X . If F_X is continuously integrable for all $x \in \mathbb{R}$ then X has a probability density f_X defined by:

$$F_X(x) = \int_{-\infty}^x f_X(y) dy \quad (65.34)$$

or alternatively:

$$f_X(x) = \lim_{\epsilon \rightarrow 0} \frac{\mathbb{P}(x \leq X \leq x + \epsilon)}{\epsilon} \quad (65.35)$$

Corollary 65.24 $\mathbb{P}(X = b) = 0, \quad \forall b \in \mathbb{R}$:

$$\mathbb{P}(X = b) = \lim_{\epsilon \rightarrow 0} \mathbb{P}(b - \epsilon < X \leq b) = \lim_{\epsilon \rightarrow 0} \int_{b-\epsilon}^b f(x) dx = 0 \quad (65.36)$$

Corollary 65.25 : From [cor. 65.24] it follows that the exact borders are not necessary:

$$\begin{aligned} \mathbb{P}(a < X < b) &= \mathbb{P}(a \leq X < b) \\ &= \mathbb{P}(a < X \leq b) = \mathbb{P}(a \leq X < \infty) \end{aligned}$$

Corollary 65.26 :

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad (65.37)$$

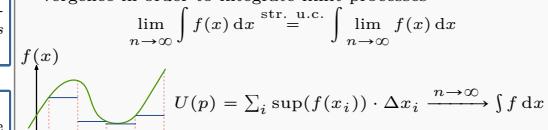
Notes

- Often the cumulative distribution function is referred to as "cdf" or simply *distribution function*.
- Often the probability density function is referred to as "pdf" or simply *density*.

5. Lebesgue Integration

Problems of Riemann Integration

- Difficult to extend to higher dimensions – general domains of definitions $f : \Omega \mapsto \mathbb{R}$
- Depends on continuity
- Integration of limit processes require strong uniform convergence in order to integrate limit processes

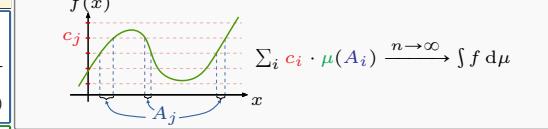


Idea

Partition domain by function values of equal size i.e. values that lie within the same sets/have the same value A_j build up the partitions w.r.t. to the variable x .

Problem: we do not know how big those sets/partitions on the x -axis will be.

Solution: we can use the measure μ of our measure space $(\Omega, \mathcal{A}, \mu)$ in order to obtain the size of our sets A_j \Rightarrow we do not have to care anymore about discontinuities, as we can measure the size of our sets using our measure.



Definition 65.29 Lebesgue Integral:

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n c_i \mu(A_i) = \int_{\Omega} f d \mu \quad \forall x \in A_i \quad (65.38)$$

Definition 65.30

Simple Functions (Random Variables): A r.v. X is called simple if it takes on only a finite number of values and hence can be written in the form:

$$X = \sum_{i=1}^n a_i \mathbb{1}_{A_i} \quad a_i \in \mathbb{R} \quad A \ni A_i = \begin{cases} 1 & \text{if } \{X = a_i\} \\ 0 & \text{else} \end{cases} \quad (65.39)$$

6. Independent Random Variables

We have seen that two events A and B are independent if knowledge that B has occurred does not change the probability that A will occur theorem 65.1.

For two random variables X, Y we want to know if knowledge of Y leaves the probability of X , to take on certain values unchanged.

Definition 65.31 Independent Random Variables:

Two real valued random variables X and Y are said to be independent iff:

$$\mathbb{P}(X \leq x | Y \leq y) = \mathbb{P}(X \leq x) \quad \forall x, y \in \mathbb{R} \quad (65.40)$$

which amounts to:

$$\begin{aligned} \mathbb{P}_{X,Y}(x, y) &= \mathbb{P}(\{X \leq x\} \cap \{Y \leq y\}) = \mathbb{P}(X \leq x, Y \leq y) \\ &= F_X(x)F_Y(y) \quad \forall x, y \in \mathbb{R} \end{aligned} \quad (65.41)$$

or alternatively iff:

$$\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A) \mathbb{P}(Y \in B) \quad \forall A, B \in \mathcal{B} \quad (65.42)$$

Note

If the joint distribution $F_{X,Y}(x,y)$ can be factorized into two functions of x and y then X and Y are independent.

Definition 65.32

Independent Identically Distributed:

10. Product Rule

Law 65.4 Product Rule: Let X, Y be two random variables then their joint probability density function is given by:

Definition 65.34 Joint Probability Distribution:

Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random vector in \mathbb{R}^n with associated cdf $F_{\mathbf{X}}$. If $F_{\mathbf{X}}$ is continuously integrable for all $x \in \mathbb{R}^n$ then \mathbf{X} has a probability density $f_{\mathbf{X}}$ defined by:

$$f_{\mathbf{X}}(x) = \int_{-\infty}^{x_n} \cdots \int_{-\infty}^{x_1} f_{\mathbf{X}}(y_1, \dots, y_n) dy_1 dy_n \quad (65.50)$$

or alternatively:

$$f_{\mathbf{X}}(\mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{\mathbb{P}(x_1 \leq X_1 \leq x_1 + \epsilon, \dots, x_n \leq X_n \leq x_n + \epsilon)}{\epsilon} \quad (65.51)$$

1. Marginal Distribution

Definition 65.35 Marginal Distribution:

14. The Expectation

Definition 65.36 Expectation:

$$\mathbb{E}[X] = \int_{\Omega} X(\omega) d\omega = \int_{\Omega} X d\mathbb{P} \quad (65.52)$$

Corollary 65.27 Expectation of simple r.v.:

If X is a simple r.v. its expectation is given by:

$$\mathbb{E}[X] = \sum_{i=1}^n a_i \mathbb{P}(A_i) \quad (65.53)$$

1. Properties

14.1.1. Linear Operators

14.1.2. Quadratic Form

Definition 65.37

proof 65.7

Expectation of a Quadratic Form:

Let $\epsilon \in \mathbb{R}^n$ be a random vector with $\mathbb{E}[\epsilon] = \mu$ and $\mathbb{V}[\epsilon] = \Sigma$:
 $\mathbb{E}[\epsilon^T A \epsilon] = \text{tr}(A\Sigma) + \mu^T A \mu$ (65.54)

2. The Jensen Inequality

Theorem 65.6 Jensen Inequality: Let X be a random variable and g some function, then it holds:

$$\begin{aligned} g(\mathbb{E}[X]) &\leq \mathbb{E}[g(X)] & \text{if } g \text{ is convex} & \text{def. 54.24} \\ g(\mathbb{E}[X]) &\geq \mathbb{E}[g(X)] & \text{if } g \text{ is concave} & \text{def. 54.25} \end{aligned} \quad (65.55)$$

3. Law of the Unconscious Statistician

Law 65.6 Law of the Unconscious Statistician:

Let $X \in \mathcal{X}, Y \in \mathcal{Y}$ be random variables where Y is defined as:
 $\mathcal{Y} = \{y | y = g(x), \forall x \in \mathcal{X}\}$

then the expectation of Y can be calculated in terms of X :

$$\mathbb{E}_Y[Y] = \mathbb{E}_X[g(x)] \quad (65.56)$$

Consequence

Hence if we \mathbb{P}_X we do not have to first calculate \mathbb{P}_Y in order to calculate $\mathbb{E}_Y[Y]$.

4. Properties

5. Law of Iterated Expectation (LIE)

Law 65.7

proof 65.8

Law of Iterated Expectation (LIE):

$$\mathbb{E}[X] = \mathbb{E}_Y \mathbb{E}[X|Y] \quad (65.57)$$

6. Hoeffding's Bound

Definition 65.38 Hoeffding's Bound:

Let $\mathbf{X} = \{X_i\}_{i=1}^n$ be i.i.d. random variables strictly bounded by the interval $[a, b]$ then it holds:

$$\mathbb{P}(|\mu_{\mathbf{X}} - \mathbb{E}[X]| \geq \epsilon) \leq 2 \exp\left(\frac{-2n^2 \epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \stackrel{[0,1]}{=} 2e^{-2n\epsilon^2} \quad (65.58)$$

Explanation 65.4. The difference of the expectation from the empirical average to be bigger than ϵ is upper bound in probability.

15. Moment Generating Function (MGF)

Definition 65.39 Moment of Random Variable: The i -th moment of a random variable X is defined as (if it exists):

$$m_i := \mathbb{E}[X^i] \quad (65.59)$$

Note

A monotonic function is required in order to satisfy inevitability.

Probability Distributions on \mathbb{R}^n

13. Joint Distribution

Definition 65.33

Joint (Cumulative) Distribution Function $F_{\mathbf{X}}$:
 Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random vector in \mathbb{R}^n , then its cumulative distribution function is defined as:

$$\begin{aligned} F_{\mathbf{X}}(\mathbf{x}) &= \mathbb{P}^X((-\infty, \mathbf{x}]) = \mathbb{P}(\mathbf{X} \leq \mathbf{x}) \\ &= \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n) \end{aligned} \quad (65.49)$$

Definition 65.40

Moment Generating Function (MGF):

$$\psi_{\mathbf{X}}(t) = \mathbb{E}[e^{tX}] \quad t \in \mathbb{R} \quad (65.60)$$

Corollary 65.28 Sum of MGF: The moment generating function of a sum of n independent variables $(X_j)_{1 \leq j \leq n}$ is the product of the moment generating functions of the components:

$$\psi_{S_n}(t) = \psi_{X_1}(t) \cdots \psi_{X_n}(t) \quad S_n := X_1 + \dots + X_n \quad (65.61)$$

Corollary 65.29 : The i -th moment of a random variable is the i -th derivative of its associated moment generating function evaluated zero:

$$\mathbb{E}[X^i] = \psi_X^{(i)}(0) \quad (65.62)$$

16. The Characteristic Function

Transforming probability distributions using the Fourier transform is a handy tool in probability in order to obtain properties or solve problems in another space before transforming them back.

Definition 65.41

Fourier Transformed Probability Measure:

$$\hat{\mu} = \int e^{i\langle u, x \rangle} \mu(dx) \quad (65.63)$$

Corollary 65.30 : As $e^{i\langle u, x \rangle}$ can be rewritten using formulaeqs. (50.9) and (50.10) it follows:

$$\hat{\mu} = \int \cos(\langle u, x \rangle) \mu(dx) + i \int \sin(\langle u, x \rangle) \mu(dx) \quad (65.64)$$

where $x \mapsto \cos(\langle x, u \rangle)$ and $x \mapsto \sin(\langle x, u \rangle)$ are both bounded and Borel i.e. Lebesgue integrable.

Definition 65.42 Characteristic Function $\varphi_{\mathbf{X}}$: Let \mathbf{X} be an \mathbb{R}^n -valued random variable. Its characteristic function $\varphi_{\mathbf{X}}$ is defined on \mathbb{R}^n as:

$$\varphi_{\mathbf{X}}(u) = \int e^{i\langle u, x \rangle} \mathbb{P}^X(dx) = \widehat{\mathbb{P}^X}(u) \quad (65.65)$$

$$= \mathbb{E}[e^{i\langle u, X \rangle}] \quad (65.66)$$

Corollary 65.31 : The characteristic function $\varphi_{\mathbf{X}}$ of a distribution always exists as it is equal to the Fourier transform of the probability measure, which always exists.

Note

This is an advantage over the moment generating function.

Theorem 65.7 : Let μ be a probability measure on \mathbb{R}^n . Then $\hat{\mu}$ is a bounded continuous function with $\hat{\mu}(0) = 1$.

Theorem 65.8 Uniqueness Theorem: The Fourier Transform $\hat{\mu}$ of a probability measure μ on \mathbb{R}^n characterizes μ . That is, if two probability measures on \mathbb{R}^n admit the same Fourier transform, they are equal.

Corollary 65.32 : Let $\mathbf{X} = (X_1, \dots, X_n)$ be an \mathbb{R}^n -valued random variable. Then the real valued r.v.'s $(X_j)_{1 \leq j \leq n}$ are independent if and only if:

$$\varphi_{\mathbf{X}}(u_1, \dots, u_n) = \prod_{j=1}^n \varphi_{X_j}(u_j) \quad (65.67)$$

Proofs

Proof 65.1: [cor. 65.11]: Let \mathcal{C} denote all open intervals. Since every open set in \mathbb{R} is the countable union of open intervals^[def. 50.12], it holds that $\sigma(\mathcal{C})$ is the Borel σ -algebra of \mathbb{R} .

Let \mathcal{D} denote all intervals of the form $(-\infty, a]$, $a \in \mathbb{Q}$.

- $(a_n)_{n>1}$ be a sequence of rationals decreasing to a and
- $(b_n)_{n>1}$ be a sequence of rationals increasing strictly to b

$(a, b) = \cup_{n=1}^{\infty} (a_n, b_n) = \cup_{n=1}^{\infty} (-\infty, b_n] \cap (-\infty, a_n]^C$

Thus $\mathcal{C} \subset \sigma(\mathcal{D})$, whence $\sigma(\mathcal{C}) \subset \sigma(\mathcal{D})$ but as each element of \mathcal{D} is a closed subset, $\sigma(\mathcal{D})$ must also be contained in the Borel sets \mathcal{B} with

$$\mathcal{B} = \sigma(\mathcal{C}) \subset \sigma(\mathcal{D}) \subset \mathcal{B}$$

Proof 65.2: theorem 65.3 Plug eq. (65.22) into the denominator and eq. (55.2) into the nominator and then use^[def. 65.18]:

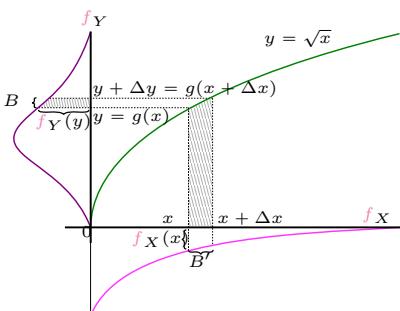
$$\frac{\mathbb{P}(B|A_j)\mathbb{P}(A_j)}{\sum_{i \in I} \mathbb{P}(B|A_i)\mathbb{P}(A_i)} = \frac{\mathbb{P}(B \cap A_j)}{\mathbb{P}(B)} = \mathbb{P}(A_j|B)$$

Proof 65.3: ??:

$$Y = g(X) \iff \mathbb{P}(Y = y) = \mathbb{P}(x \in \mathcal{X}_y) = \mathbb{P}_Y(y)$$

Proof 65.4: ?? (non-formal): The probability contained in a differential area must be invariant under a change of variables that is:

$$|\mathbb{f}_Y(y) dy| = |\mathbb{f}_X(x) dx|$$



Proof 65.5: ?? from CDF:

$$\mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \begin{cases} \mathbb{P}(X \leq g^{-1}(y)) & \text{if } g \text{ is increasing} \\ \mathbb{P}(X \geq g^{-1}(y)) & \text{if } g \text{ is decreasing} \end{cases}$$

If g is monotonically increasing:

$$F_Y(y) = F_X(g^{-1}(y))$$

$$\mathbb{f}_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = \mathbb{f}_X(x) \cdot \frac{d}{dy} g^{-1}(y)$$

If g is monotonically decreasing:

$$F_Y(y) = 1 - F_X(g^{-1}(y))$$

$$\mathbb{f}_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = -\mathbb{f}_X(x) \cdot \frac{d}{dy} g^{-1}(y)$$

Proof 65.6: ?? Let $B = [x, x + \Delta x]$ and $B' = [y, y + \Delta y] = [g(x), g(x + \Delta x)]$ we know that the probability of equal events is equal:

$$y = g(x) \Rightarrow \mathbb{P}(y) = \mathbb{P}(g(x)) \text{ (for disc. rv.)}$$

Now lets consider the probability for the continuous r.v.s:

$$\mathbb{P}(X \in B) = \int_x^{x+\Delta x} f_X(t) dt \xrightarrow{\Delta x \rightarrow 0} |\Delta x \cdot f_X(x)|$$

For y we use Taylor (??)

$$\begin{aligned} g(x + \Delta x) &\stackrel{\text{eq. (54.56)}}{=} g(x) + \frac{dg}{dx} \Delta y \quad \text{for } \Delta x \rightarrow 0 \\ &= y + \Delta y \quad \text{with } \Delta y := \frac{dg}{dx} \cdot \Delta x \end{aligned} \quad (65.68)$$

Thus for $\mathbb{P}(Y \in B')$ it follows:

$$\begin{aligned} \mathbb{P}(X \in B') &= \int_y^{y+\Delta y} f_Y(t) dt \xrightarrow{\Delta y \rightarrow 0} |\Delta y \cdot f_Y(y)| \\ &= \left| \frac{dg}{dx}(x) \Delta x \cdot f_Y(y) \right| \end{aligned}$$

Now we simply need to related the surface of the two pdfs:

$$\begin{aligned} B &= [x, x + \Delta x] \underset{\text{same surfaces}}{\propto} [y, y + \Delta y] = B' \\ \mathbb{P}(Y \in B) &= \mathbb{P}(X \in B') \\ \xrightarrow{\Delta y \rightarrow 0} |f_Y(y) \cdot \Delta y| &= \left| f_Y(y) \cdot \frac{dg}{dx}(x) \Delta x \right| = |f_X(x) \cdot \Delta x| \\ f_Y(y) \left| \frac{dg}{dx}(x) \right| |\Delta x| &= f_X(x) \cdot |\Delta x| \\ \Rightarrow f_Y(y) &= \frac{f_X(x)}{\left| \frac{dg}{dx}(x) \right|} = \frac{f_X(g^{-1}(y))}{\left| \frac{dg}{dx}(x) \right|} \end{aligned}$$

Proof 65.7: [def. 65.37]

$$\begin{aligned} \mathbb{E}[\epsilon^T A \epsilon] &\stackrel{\text{eq. (59.54)}}{=} \mathbb{E}[\text{tr}(\epsilon^T A \epsilon)] \\ &\stackrel{\text{eq. (59.56)}}{=} \mathbb{E}[\text{tr}(A \epsilon \epsilon^T)] \\ &= \text{tr}(\mathbb{E}[A \epsilon \epsilon^T]) \\ &= \text{tr}(A \mathbb{E}[\epsilon \epsilon^T]) \\ &= \text{tr}(A(\Sigma + \mu \mu^T)) \\ &= \text{tr}(A\Sigma) + \text{tr}(A\mu \mu^T) \\ &\stackrel{\text{eq. (59.54)}}{=} \text{tr}(A\Sigma) + A\mu \mu^T \end{aligned}$$

Proof 65.8: law 65.7

$$\begin{aligned} \mathbb{E}[X] &= \sum_x x \cdot \mathbb{P}_X(x) = \sum_x x \cdot \sum_{y \in \mathcal{Y}} \mathbb{P}_{X,Y}(x,y) \\ &= \sum_x x \cdot \sum_y \mathbb{P}_{X|Y}(x|y) \cdot \mathbb{P}_Y(y) \\ &= \sum_y \mathbb{P}_Y(y) \cdot \sum_x x \cdot \mathbb{P}_{X|Y}(x|y) \\ &= \sum_y \mathbb{P}_Y(y) \cdot \mathbb{E}[X|Y] = \mathbb{E}_Y[\mathbb{E}[X|Y]] \end{aligned}$$

Examples

- Example 65.1 :**
- Toss of a coin (with head and tail): $\Omega = \{H, T\}$.
 - Two tosses of a coin: $\Omega = \{HH, HT, TH, TT\}$
 - A cubic die: $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6\}$
 - The positive integers: $\Omega = \{1, 2, 3, \dots\}$
 - The reals: $\Omega = \{\omega | \omega \in \mathbb{R}\}$

Example 65.2 :

- Head in coin toss $A = \{H\}$
- Odd number in die roll: $A = \{\omega_1, \omega_3, \omega_5, \dots\}$
- The integers smaller five: $A = \{1, 2, 3, 4\}$

Example 65.3 : If the sample space is a die toss $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6\}$, the sample space may be that we are only told whether an even or odd number has been rolled:

$$\mathcal{F} = \{\emptyset, \{\omega_1, \omega_3, \omega_5\}, \{\omega_2, \omega_4, \omega_6\}\}$$

Example 65.4 : If we are only interested in the subset $A \in \Omega$ of our experiment, then we can look at the corresponding generating σ -algebra $\sigma(A) = \{\emptyset, A, A^C, \Omega\}$.

Example 65.5 :

- open half-lines: $(-\infty, a)$ and (a, ∞) ,
- union of open half-lines: $(a, b) = (-\infty, a) \cup (b, \infty)$,
- closed interval: $[a, b] = \overline{(-\infty, \cup a) \cup (b, \infty)}$,
- closed half-lines:

 - $(-\infty, a] = \bigcup_{n=1}^{\infty} [a - n, a]$ and $[a, \infty) = \bigcup_{n=1}^{\infty} [a, a + n]$,

- half-open and half-closed $[a, b] = (-\infty, b] \cup (a, \infty)$,
- every set containing only one real number: $\{a\} = \bigcap_{n=1}^{\infty} (a - \frac{1}{n}, a + \frac{1}{n})$,
- every set containing finitely many real numbers: $\{a_1, \dots, a_n\} = \bigcup_{k=1}^n a_k$.

Note: why do we need probability density functions

A continuous random variable X can realise an infinite count of real number values within its support B (as there are an infinitude of points in a line segment).

Thus we have an infinitude of values whose sum of probabilities must equal one.

Thus these probabilities must each be zero otherwise we would obtain probability of ∞ . As we can not work with zero probabilities we use the next best thing, infinitesimal probabilities (defined as a limit).

We say they are almost surely equal to zero:

$$\mathbb{P}(X = x) = 0 \quad \text{a.s.}$$

To have a sensible measure of the magnitude of these infinitesimal quantities, we use the concept of probability density, which yields a probability mass when integrated over an interval.

Example 65.6 Equivalent (Probability) Measures:

$$\begin{aligned} \Omega &= \{1, 2, 3\} & \mathbb{P}(\{1, 2, 3\}) &= \{2/3, 1/6, 1/6\} \\ & & \tilde{\mathbb{P}}(\{1, 2, 3\}) &= \{1/3, 1/3, 1/3\} \end{aligned}$$

Example 65.7 :

Example 65.8 ??: Let $X, Y \sim \mathcal{N}(0, 1)$.

Question: proof that:

$$U = X + Y \quad V = X - 1$$

are independent and normally distributed:

$$h(u, v) = \begin{cases} h_1(u, v) = \frac{u+v}{\sqrt{2}} \\ h_2(u, v) = \frac{u-v}{\sqrt{2}} \end{cases} \quad J = \det \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{bmatrix} = -\frac{1}{2}$$

$$\begin{aligned} f_{U,V} &= f_{X,Y}(x, y) \cdot \frac{1}{2} \\ &\stackrel{\text{indp.}}{=} f_X(x) \cdot f_Y(y) \\ &= \frac{1}{2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \\ &= \frac{1}{2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\left(\frac{u+v}{\sqrt{2}}\right)^2} \cdot \frac{1}{2} e^{-\left(\frac{u-v}{\sqrt{2}}\right)^2} \\ &= \frac{1}{\sqrt{2\pi}\sqrt{2}} e^{-\frac{u^2}{4}} \cdot \frac{1}{\sqrt{2\pi}\sqrt{2}} e^{-\frac{v^2}{4}} \end{aligned}$$

Thus U, V are independent r.v. distributed as $\mathcal{N}(0, 2)$.

Statistics

The probability that a discret random variable x is equal to some value $\bar{x} \in \mathcal{X}$ is:

$$p_x(\bar{x}) = \mathbb{P}(x = \bar{x})$$

Definition 66.1 Almost Surely \mathbb{P} -(a.s.):

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. An event $\omega \in \mathcal{F}$ happens almost surely iff

$$\mathbb{P}(\omega) = 1 \iff \omega \text{ happens a.s.} \quad (66.1)$$

Definition 66.2 Probability Mass Function (PMF):

Definition 66.3 Discrete Random Variable (DVR): The set of possible values \bar{x} of \mathcal{X} is countable of finite. $\mathcal{X} = \{0, 1, 2, 3, 4, \dots, 8\} \quad \mathcal{X} = \mathbb{N} \quad (66.2)$

Definition 66.4 Probability Density Function (PDF): Is real function $f: \mathbb{R}^n \rightarrow [0, \infty)$ that satisfies:

Non-negativity: $f(x) \geq 0, \quad \forall x \in \mathbb{R}^n \quad (66.3)$

Normalization: $\int_{-\infty}^{\infty} f(x) dx = 1 \quad (66.4)$

Must be integrable (66.5)

Property 66.3: **Monotonically Increasing**

$$x \leq y \iff F_X(x) \leq F_X(y) \quad \forall x, y \in \mathbb{R} \quad (66.10)$$

Upper Limit

$$\lim_{x \rightarrow \infty} F_X(x) = 1 \quad (66.11)$$

Lower Limit

$$\lim_{x \rightarrow -\infty} F_X(x) = 0 \quad (66.12)$$

Definition 66.8 CDF of a discret rv X : Let X be discret rv with pdf \mathbb{P}_X , then the CDF of X is given by:

$$F_X(x) = \mathbb{P}(X \leq x) = \sum_{t=-\infty}^x \mathbb{P}_X(t) \quad (66.13)$$

Definition 66.9 CDF of a continuous rv X : Let X be continuous rv with pdf f_X , then the CDF of X is given by:

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \iff \frac{\partial F_X(x)}{\partial x} = f_X(x) \quad (66.14)$$

Lemma 66.1 Probability Interval: Let X be a continuous rrv with pdf f_X and cumulative distribution function F_X , then it holds that:

$$\mathbb{P}(a \leq X \leq b) = F_X(b) - F_X(a) \quad (66.13)$$

Proof 66.3: [def. 66.9]:

$$F_X(x) = \mathbb{P}(X \leq x) = \mathbb{P}(X \in (-\infty, x)) = \int_{-\infty}^x f_X(t) dt$$

Proof 66.4: lemma 66.1:

$$\mathbb{P}(a \leq X \leq b) = \mathbb{P}(X \leq b) - \mathbb{P}(X \leq a)$$

or by the fundamental theorem of calculus (theorem 54.2):

$$\mathbb{P}(a \leq X \leq b) = \int_a^b f_X(t) dt = \int_a^b \frac{\partial F_X(t)}{\partial t} dt = [F_X(t)]_a^b$$

Theorem 66.2 A continuous rv is fully characterized by its CDF: A cumulative distribution function completely determines the distribution of a continuous real-valued random variable.

1. Key figures

1. The Expectation

Definition 66.10 Expectation (disc. case): $\mathbb{E}_X := \mathbb{E}_x[x] := \sum_{x \in \mathcal{X}} \bar{x} p_x(\bar{x}) \quad (66.14)$

Definition 66.11 Expectation (cont. case): $\mathbb{E}_x[x] := \int_{\mathcal{X}} \bar{x} f_x(\bar{x}) d\bar{x} \quad (66.15)$

Law 66.1 Expectation of independent variables:

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \quad (66.16)$$

Property 66.4 Translation and scaling: If $\mathbf{X} \in \mathbb{R}^n$ and $\mathbf{Y} \in \mathbb{R}^m$ are random vectors, and $a, b, \mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ are constants then it holds:

$$\mathbb{E}[a + b\mathbf{X} + c\mathbf{Y}] = a + b\mathbb{E}[\mathbf{X}] + c\mathbb{E}[\mathbf{Y}] \quad (66.17)$$

Thus \mathbb{E} is a linear operator ([def. 54.15]).

Note: Expectation of the expectation

The expectation of a r.v. X is a constant hence with Property 66.6 it follows:

$$\mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X] \quad (66.18)$$

Property 66.5 Matrix \times Expectation: If $\mathbf{X} \in \mathbb{R}^n$ is a random vector and $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times m}$ are constant matrices then it holds:

$$\mathbb{E}[\mathbf{AXB}] = \mathbf{A}\mathbb{E}[(\mathbf{XB})] = \mathbf{A}\mathbb{E}[\mathbf{X}]\mathbf{B} \quad (66.19)$$

Definition 66.7 Cumulative distribution function (CDF): Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space.

The (cumulative) distribution function of a real-valued random variable X is the function given by:

$$F_X(x) = \mathbb{P}(X \leq x) \quad \forall x \in \mathbb{R}$$

Proof 66.5: eq. (66.24):

$$\begin{aligned}\mathbb{E}[XY] &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y)xy \\ ?? &= \sum_{x \in \mathcal{X}} p_X(x)x \sum_{y \in \mathcal{Y}} p_Y(y)y = \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

Definition 66.12

Autocorrelation/Crosscorelation $\gamma(t_1, t_2)$: Describes the covariance (def. 66.16) between the two values of a stochastic process $(\mathbf{X}_t)_{t \in T}$ at different time points t_1 and t_2 .

$$\gamma(t_1, t_2) = \text{Cov}[\mathbf{X}_{t_1}, \mathbf{X}_{t_2}] = \mathbb{E}[(\mathbf{X}_{t_1} - \mu_{t_1})(\mathbf{X}_{t_2} - \mu_{t_2})] \quad (66.20)$$

For zero time differences $t_1 = t_2$ the autocorrelation function equals the variance:

$$\gamma(t, t) = \text{Cov}[\mathbf{X}_t, \mathbf{X}_t] \stackrel{\text{eq. (66.35)}}{=} \mathbb{V}[\mathbf{X}_t] \quad (66.21)$$

Notes

- Hence the autocorrelation describes the correlation of a function or signal with itself at a previous time point.
- Given a random time dependent variable $\mathbf{x}(t)$ the autocorrelation function $\gamma(t, t - \tau)$ describes how similar the time translated function $\mathbf{x}(t - \tau)$ and the original function $\mathbf{x}(t)$ are.
- If there exists some relation between the values of the time series that is non-random then the autocorrelation is non-zero.
- The autocorrelation is maximized/most similar for no translation $\tau = 0$ at all.

2. Key Figures

1. The Expectation

Definition 66.13 Expectation (disc. case):

$$\mu_X := \mathbb{E}_x[x] := \sum_{\mathbf{x} \in \mathcal{X}} \bar{x} p_x(\bar{x}) \quad (66.22)$$

Definition 66.14 Expectation (cont. case):

$$\mathbb{E}_x[x] := \int_{\mathcal{X} \in \mathcal{X}} \bar{x} f_x(\bar{x}) d\bar{x} \quad (66.23)$$

Law 66.2 Expectation of independent variables:

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \quad (66.24)$$

Property 66.6 Translation and scaling: If $\mathbf{X} \in \mathbb{R}^n$ and $\mathbf{Y} \in \mathbb{R}^m$ are random vectors, and $a, b, \mathbf{a} \in \mathbb{R}^n$ are constants then it holds:

$$\mathbb{E}[a + b\mathbf{X} + c\mathbf{Y}] = a + b\mathbb{E}[\mathbf{X}] + c\mathbb{E}[\mathbf{Y}] \quad (66.25)$$

Thus \mathbb{E} is a linear operator (def. 54.15).

Property 66.7

Affine Transformation of the Expectation:

If $\mathbf{X} \in \mathbb{R}^n$ is a random vector, $\mathbf{A} \in \mathbb{R}^{m \times n}$ a constant matrix and $b \in \mathbb{R}^m$ then it holds:

$$\mathbb{E}[\mathbf{AX} + b] = \mathbf{A}\mu + b \quad (66.26)$$

Note: Expectation of the expectation

The expectation of a r.v. X is a constant hence with Property 66.6 it follows:

$$\mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X] \quad (66.27)$$

Property 66.8 Matrix×Expectation: If $\mathbf{X} \in \mathbb{R}^n$ is a random vector and $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times m}$ are constant matrices then it holds:

$$\mathbb{E}[\mathbf{AXB}] = \mathbf{A}\mathbb{E}[(\mathbf{XB})] = \mathbf{A}\mathbb{E}[\mathbf{X}]\mathbf{B} \quad (66.28)$$

Proof 66.6: eq. (66.24):

$$\begin{aligned}\mathbb{E}[XY] &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y)xy \\ ?? &= \sum_{x \in \mathcal{X}} p_X(x)x \sum_{y \in \mathcal{Y}} p_Y(y)y = \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

2. The Variance

Definition 66.15 Variance $\mathbb{V}[X]$: The variance of a random variable X is the expected value of the squared deviation from the expectation of X ($\mu = \mathbb{E}[X]$).

It is a measure of how much the actual values of a random variable X fluctuate around its expected value $\mathbb{E}[X]$ and is defined by:

$$\mathbb{V}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2] \stackrel{\text{see ?? 66.7}}{=} \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (66.29)$$

2.2.1. Properties

Property 66.9 Variance of a Constant: If $a \in \mathbb{R}$ is a constant then it follows that its expected value is deterministic \Rightarrow we have no uncertainty \Rightarrow no variance:

$$\mathbb{V}[a] = 0 \quad \text{with} \quad a \in \mathbb{R} \quad (66.30)$$

see shift and scaling for proof ?? 66.8

Property 66.10 Shifting and Scaling:

$$\mathbb{V}[a + bX] = a^2\sigma^2 \quad \text{with} \quad a \in \mathbb{R} \quad (66.31)$$

see ?? 66.8

Property 66.11 [proof 66.9]

Affine Transformation of the Variance: If $\mathbf{X} \in \mathbb{R}^n$ is a random vector, $\mathbf{A} \in \mathbb{R}^{m \times n}$ a constant matrix and $b \in \mathbb{R}^m$ then it holds:

$$\mathbb{V}[\mathbf{AX} + b] = \mathbf{A}\mathbb{V}[\mathbf{X}]\mathbf{A}^\top \quad (66.32)$$

Definition 66.16 Covariance: The Covariance is a measure of how much two or more random variables vary linearly with each other.

$$\begin{aligned}\text{Cov}[\mathbf{X}, \mathbf{Y}] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]\end{aligned} \quad (66.33)$$

see ?? 66.10

Definition 66.17 Covariance Matrix: The variance of a k -dimensional random vector $\mathbf{X} = (X_1 \dots X_k)$ is given by a p.s.d. eq. (59.109) matrix called Covariance Matrix. The Covariance is a measure of how much two or more random variables vary linearly with each other and the Variance on the diagonal is again a measure of how much a variable varies:

$$\begin{aligned}\mathbb{V}[\mathbf{X}] &:= \Sigma(\mathbf{X}) := \text{Cov}[\mathbf{X}, \mathbf{X}] := \\ &= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \\ &= \mathbb{E}[\mathbf{XX}^\top] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^\top \in [-\infty, \infty]\end{aligned} \quad (66.34)$$

$$\begin{aligned}&= \begin{bmatrix} \mathbb{V}[X_1] & \cdots & \mathbb{C}ov[X_1, X_k] \\ \vdots & \ddots & \vdots \\ \mathbb{C}ov[X_k, X_1] & \cdots & \mathbb{V}[X_k] \end{bmatrix} \\ &= \begin{bmatrix} \mathbb{E}[(X_1 - \mu_1)(X_1 - \mu_1)] & \cdots & \mathbb{E}[(X_1 - \mu_1)(X_k - \mu_k)] \\ \vdots & \ddots & \vdots \\ \mathbb{E}[(X_k - \mu_k)(X_1 - \mu_1)] & \cdots & \mathbb{E}[(X_k - \mu_k)(X_k - \mu_k)] \end{bmatrix}\end{aligned}$$

Note: Covariance and Variance

The variance is a special case of the covariance in which two variables are identical:

$$\text{Cov}[\mathbf{X}, \mathbf{X}] = \mathbb{V}[\mathbf{X}] = \sigma^2(X) \equiv \sigma_X^2 \quad (66.35)$$

Property 66.12 Translation and Scaling:

$$\text{Cov}(a + bX, c + dY) = bd\text{Cov}(X, Y) \quad (66.36)$$

Property 66.13

Affine Transformation of the Covariance:

If $\mathbf{X} \in \mathbb{R}^n$ is a random vector, $\mathbf{A} \in \mathbb{R}^{m \times n}$ a constant matrix and $b \in \mathbb{R}^m$ then it holds:

$$\text{Cov}[\mathbf{AX} + b] = \mathbf{A}\mathbb{V}[\mathbf{X}]\mathbf{A}^\top = \mathbf{A}\Sigma(\mathbf{X})\mathbf{A}^\top \quad (66.37)$$

Definition 66.18 Correlation Coefficient: Is the standardized version of the covariance:

$$\text{Corr}[\mathbf{X}] := \frac{\text{Cov}[\mathbf{X}]}{\sigma_{X_1} \cdots \sigma_{X_k}} \in [-1, 1] \quad (66.38)$$

$$= \begin{cases} +1 & \text{if } Y = aX + b \text{ with } a > 0, b \in \mathbb{R} \\ -1 & \text{if } Y = aX + b \text{ with } a < 0, b \in \mathbb{R} \end{cases}$$

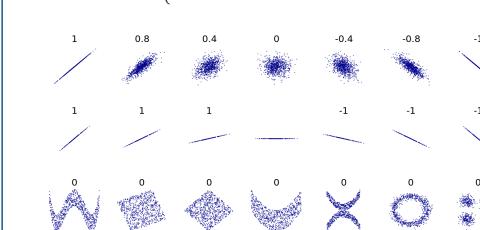


Figure 17: Several sets of (x, y) points, with their correlation coefficient

Law 66.3 Translation and Scaling:

$$\text{Corr}(a + bX, c + dY) = \text{sign}(b)\text{sign}(d)\text{Cov}(X, Y) \quad (66.39)$$

Note

- The correlation/covariance reflects the noisiness and direction of a linear relationship (top row fig. 17), but not the slope of that relationship (middle row fig. 17) nor many aspects of nonlinear relationships (bottom row).
- The set in the center of fig. 17 has a slope of 0 but in that case the correlation coefficient is undefined because the variance of Y is zero.
- Zero covariance/correlation $\text{Cov}(X, Y) = \text{Corr}(X, Y) = 0$ implies that there does not exist a linear relationship between the random variables X and Y .

Difference Covariance&Correlation

- Variance is affected by scaling and covariance not ?? and law 66.3.
- Correlation is dimensionless, whereas the unit of the covariance is obtained by the product of the units of the two RV variables.

Law 66.4 Covariance of independent RVs: The covariance/correlation of two independent variable's ?? is zero:

$$\begin{aligned}\text{Cov}[X, Y] &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \\ \stackrel{\text{eq. (66.24)}}{=} &= \mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[X]\mathbb{E}[Y] = 0\end{aligned}$$

Zero covariance/correlation \Rightarrow independence

$$\text{Cov}(X, Y) = \text{Corr}(X, Y) = 0 \Rightarrow p_{X,Y}(x, y) = p_X(x)p_Y(y)$$

For example: let $X \sim \mathcal{U}([-1, 1])$ and let $Y = X^2$.

1. Clearly X and Y are dependent

2. But the covariance/correlation between X and Y is non-zero:
 $\text{Cov}(X, Y) = \text{Cov}(X, X^2) = \mathbb{E}[X \cdot X^2] - \mathbb{E}[X]\mathbb{E}[X^2]$
 $= \mathbb{E}[X^3] - \mathbb{E}[X]\mathbb{E}[X^2] \stackrel{\text{eq. (66.63)}}{=} 0 - 0 \cdot \mathbb{E}[X^2] \stackrel{\text{eq. (66.52)}}{=} 0 \Rightarrow$
 \Rightarrow the relationship between Y and X must be non-linear.

Definition 66.19 Quantile: Are specific values q_α in the range (def. 54.10) of a random variable X that are defined as the value for which the cumulative probability is less than q_α with probability $\alpha \in (0, 1)$:

$$q_\alpha : \mathbb{P}(X \leq x) = F_X(q_\alpha) = \alpha \xrightarrow{\text{F invert.}} q_\alpha = F_X^{-1}(\alpha) \quad (66.40)$$

3. Proofs

Proof 66.7: eq. (66.29)

$$\mathbb{V}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2]$$

$$\stackrel{\text{Property 66.6}}{=} \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 = \mathbb{E}[X^2] - \mu^2$$

Proof 66.8: Property 66.10

$$\begin{aligned}\mathbb{V}[a + bX] &= \mathbb{E}[(a + bX - \mathbb{E}[a + bX])^2] \\ &= \mathbb{E}\left[\left(\cancel{a} + \cancel{b}X - \cancel{a} - \cancel{b}\mathbb{E}[X]\right)^2\right] \\ &= \mathbb{E}[(bX - b\mathbb{E}[X])^2] \\ &= \mathbb{E}[b^2(X - \mathbb{E}[X])^2] \\ &= b^2\mathbb{E}[(X - \mathbb{E}[X])^2] = b^2\sigma^2\end{aligned}$$

Proof 66.9: Property 66.11

$$\begin{aligned}\text{Cov}[\mathbf{AX} + b] &= \mathbb{E}[(\mathbf{AX} + b - \mathbb{E}[\mathbf{AX} + b])^2] + 0 = \\ &= \mathbb{E}[(\mathbf{AX} - \mathbb{E}[\mathbf{AX}])(\mathbf{AX} - \mathbb{E}[\mathbf{AX}])^\top] \\ &= \mathbb{E}[\mathbf{A}(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{A}(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top)] \\ &= \mathbb{E}[\mathbf{A}(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - (\mathbb{E}[\mathbf{X}])^\top \mathbf{A}^\top)] \\ &= \mathbb{E}[\mathbf{A}((\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - (\mathbb{E}[\mathbf{X}])^\top)^\top \mathbf{A}^\top] = \mathbf{A}\mathbb{V}[\mathbf{X}]\mathbf{A}^\top\end{aligned}$$

Proof 66.10: eq. (66.33)

$$\begin{aligned}\text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] + \mathbb{E}[X]\mathbb{E}[Y] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[X]\mathbb{E}[Y] + \mathbb{E}[X]\mathbb{E}[Y] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

Discrete Distributions

Definition 66.20 Multivariate Distribution: the variate refers to the number of input variables i.e. a m-variate distribution has m-input variables whereas a uni-variate distribution has only one.

Dimensional vs. Multivariate

The dimension refers to the number of dimensions we need to embed the function. If the variables of a function are independent than the dimension is the same as the number of inputs but the number of input variables can also be less.

1. Bernoulli Distribution

$\text{Bern}(p)$

Definition 66.21 Bernoulli Trial: Is a random experiment with exactly two possible outcomes, success (1) and failure (0), in which the probability of success/failure is constant in every trial i.e. independent trials.

Definition 66.22 Bernoulli Distribution $X \sim \text{Bern}(p)$:

X is a binary variable i.e. can only attain the values 0 (failure) or 1 (success) with a parameter p that signifies the success probability:

$$p(x; p) = \begin{cases} p & \text{for } x = 1 \\ 1 - p & \text{for } x = 0 \end{cases} \iff \begin{cases} P(X = 1) = p \\ P(X = 0) = 1 - p \end{cases}$$

$$= p^x \cdot (1 - p)^{1-x} \quad \text{for } x \in \{0, 1\}$$

$$\mathbb{E}[X] = p \quad (66.41) \quad \mathbb{V}[X] = p(1 - p) \quad (66.42)$$

2. Multinoulli/Categorical Distribution

$\text{Cat}(n, p)$

Definition 66.23

Multinoulli/Categorical Distribution $X \sim \text{Cat}(p)$: Is the generalization of the Bernoulli distribution^[def. 66.22] to a sample space^[def. 65.2] of k individual items $\{c_1, \dots, c_k\}$ with probabilities $p = \{p_1, \dots, p_k\}$:

$$p(x = c_i | p) = p_i \iff p(x | p) = \prod_i p_i^{\delta[x=c_i]}$$

$$\sum_{j=1}^k p_j = 1 \quad p_j \in [0, 1] \quad \forall j = 1, \dots, k \quad (66.43)$$

$$\mathbb{E}[X] = p \quad \mathbb{V}[X]_{i,j} = \Sigma_{i,j} = \begin{cases} p_i(1 - p_i) & \text{if } i = j \\ -p_i p_j & \text{if } i \neq j \end{cases}$$

Corollary 66.3

One-hot encoded Categorical Distribution: If we encode the k categories by a sparse vectors^[def. 59.70] with norm one:

$$\mathbb{B}_r^n = \left\{ \mathbf{x} \in \{0, 1\}^n : \mathbf{x}^\top \mathbf{x} = \sum_{i=1}^n \mathbf{x}_i = 1 \right\}$$

s.t. $\mathbf{x}_j = \mathbf{e}_j \iff \mathbf{x} = \mathbf{c}_j$

then we can rewrite eq. (66.43) as:

$$p(\mathbf{x} | p) = \prod_i \mathbf{x}_i \cdot p_i \quad \sum_{j=1}^k p_j = 1 \quad (66.44)$$

3. Binomial Distribution

$\text{B}(n, p)$

Definition 66.24 Binomial Coefficient:

The binomial coefficient occurs inside the binomial distribution^[def. 66.25] and signifies the different combinations/order that x out of n successes can happen see also [def. 64.5].

Definition 66.25 Binomial Distribution [proof 66.1]:

Models the probability of exactly x success given a fixed number of n -Bernoulli experiments^[def. 66.21], where the probability of success, of a single experiment is given by p :

$$p(x) = \binom{n}{x} p^x (1 - p)^{n-x} \quad \begin{matrix} n: \text{nb. of repetitions} \\ x: \text{nb. of successes} \\ p: \text{probability of success} \end{matrix}$$

$$\mathbb{E}[X] = np \quad (66.45) \quad \mathbb{V}[X] = np(1 - p) \quad (66.46)$$

Explanation 66.1. Lets consider a box of n balls consisting of black and white balls. If we want to know the probability of drawing first exactly x white and then exactly $n - x$ black balls we can simply calculate:

$$\underbrace{(p \cdots p)}_{x\text{-times}} \cdot \underbrace{(q \cdots q)}_{n-x\text{-times}} = p^x q^{n-x}$$

But if we do not care about the order we need to increase the probability by the different ways we can achieve this result by adding the binomial coefficient.

4. Geometric Distribution

$\text{Geom}(p)$

Definition 66.26 Geometric Distribution $\text{Geom}(p)$:

Models the probability of the number X of Bernoulli trials^[def. 66.21] until the first success

$$p(x) = p(1 - p)^{x-1} \quad \begin{matrix} x: \text{nb. of repetitions until first} \\ \text{success} \\ p: \text{success probability of single Bernoulli experiment} \end{matrix}$$

$$F(x) = \sum_{i=1}^x p(1 - p)^{i-1} \stackrel{\text{eq. (51.4)}}{=} 1 - (1 - p)^x$$

$$\mathbb{E}[X] = \frac{1}{p} \quad (66.47) \quad \mathbb{V}[X] = \frac{1-p}{p^2} \quad (66.48)$$

Notes

- $\mathbb{E}[X]$ is the mean waiting time until the first success
- the number of trials x in order to have at least one success with a probability of $p(x)$:

$$x \geq \frac{p(x)}{1 - p}$$

- $\log(1 - p) \approx -p$ for small p

5. Poisson Distribution

$\text{Pois}(\lambda)$

Definition 66.27 Poisson Distribution: Is an extension of the binomial distribution, where the realization x of the random variable X may attain values in $\mathbb{Z}_{\geq 0}$.

It expresses the probability of a given number of events X occurring in a fixed interval if those events occur independently of the time since the last event.

$$p(x) = e^{-\lambda} \frac{\lambda^x}{x!} \quad \lambda > 0 \quad x \in \mathbb{Z}_{\geq 0} \quad (66.49)$$

Event Rate λ : describes the average number of events in a single interval.

$$\mathbb{E}[X] = \lambda \quad (66.50) \quad \mathbb{V}[X] = \lambda \quad (66.51)$$

Continuous Distributions

1. Uniform Distribution

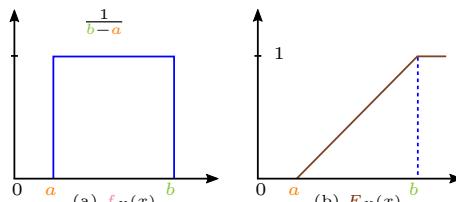
$$\mathcal{U}(a, b)$$

Definition 66.28 Uniform Distribution $\mathcal{U}(a, b)$: Is probability distribution, where all intervals of the same length on the distribution's support^[def. 66.6] $\text{supp}(\mathcal{U}[a, b]) = [a, b]$ are equally probable/likely.

$$f(x) = \frac{1}{b-a} \mathbb{1}_{x \in [a,b]} = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{else} \end{cases} \quad (66.52)$$

$$F(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & a \leq x \leq b \\ 1 & x > b \end{cases} \quad (66.53)$$

$$\mathbb{E}[X] = \frac{a+b}{2} \quad \mathbb{V}(X) = \frac{(b-a)^2}{12} \quad (66.54)$$



2. Exponential Distribution

$$\exp(\lambda)$$

Definition 66.29 Exponential Distribution $X \sim \exp(\lambda)$: Is the continuous analogue to the geometric distribution [def. 66.26].

It describes the probability $f(x; \lambda)$ that a continuous Poisson process (i.e., a process in which events occur continuously and independently at a constant average rate) will succeed/change state after a time interval x .

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \quad (66.55)$$

$$F(x; \lambda) = \begin{cases} 1 - e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \quad (66.56)$$

$$\mathbb{E}[X] = \frac{1}{\lambda} \quad \mathbb{V}(X) = \frac{1}{\lambda^2} \quad (66.57)$$

3. Laplace Distribution

Definition 66.30 Laplace Distribution:

Laplace Distribution $f(x; \mu, \sigma) = \frac{1}{2\sigma} \exp\left(-\frac{|x-\mu|}{\sigma}\right)$ (66.58)

4. The Normal Distribution

$$\mathcal{N}(\mu, \sigma)$$

Definition 66.31 Normal Distribution $X \sim \mathcal{N}(\mu, \sigma^2)$: Is a symmetric distribution where the population parameters μ, σ^2 are equal to the expectation and variance of the distribution:

$$\mathbb{E}[X] = \mu \quad \mathbb{V}(X) = \sigma^2 \quad (66.59)$$

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2\right\} \quad (66.60)$$

$$F(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left\{-\frac{1}{2} \left(\frac{u-\mu}{\sigma}\right)^2\right\} du \quad (66.61)$$

$$x \in \mathbb{R} \quad \text{or} \quad -\infty < x < \infty$$

$$\varphi_X(u) = \exp\left\{iu\mu - \frac{u^2\sigma^2}{2}\right\} \quad (66.62)$$

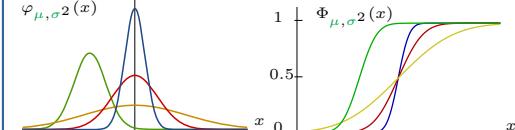


Figure 19:
 $\mu = 0 \quad \mu = 0 \quad \mu = 0 \quad \mu = -2$
 $\sigma^2 = 0.2 \quad \sigma^2 = 1.0 \quad \sigma^2 = 5.0 \quad \sigma^2 = 0.5$

Property 66.14: $P_X(\mu - \sigma \leq x \leq \mu + \sigma) = 0.66$

Property 66.15: $P_X(\mu - 2\sigma \leq x \leq \mu + 2\sigma) = 0.95$

5. The Standard Normal distribution

$$\mathcal{N}(0, 1)$$

Historic Problem: the cumulative distribution eq. (66.61) does not have an analytical solution and numerical integration was not always computationally so easy. So how should people calculate the probability of x falling into certain ranges $\mathbb{P}(x \in [a, b])$?

Solution: use a standardized form/set of parameters (by convention) $\mathcal{N}(0, 1)$ and tabulate many different values for its cumulative distribution $\Phi(x)$ s.t. we can transform all families of Normal Distributions into the standardized version $\mathcal{N}(\mu, \sigma^2) \xrightarrow{z} \mathcal{N}(0, 1)$ and look up the value in its table.

Definition 66.32

Standard Normal Distribution $X \sim \mathcal{N}(0, 1)$:

$$\mathbb{E}[X] = 0 \quad \mathbb{V}(X) = 1 \quad (66.63)$$

$$f(x; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \quad (66.64)$$

$$F(x; 0, 1) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}u^2} du \quad (66.65)$$

$$x \in \mathbb{R} \quad \text{or} \quad -\infty < x < \infty$$

$$\psi_X(u) = e^{\frac{u^2}{2}} \quad \varphi_X(u) = e^{-\frac{u^2}{2}} \quad (66.66)$$

Corollary 66.4

Standard Normal Distribution Notation: As the standard normal distribution is so commonly used people often use the letter Z in order to denote its the *standard* normal distribution and its α -quantile^[def. 66.19] is then denoted by:

$$z_\alpha = \Phi^{-1}(\alpha) \quad \alpha \in (0, 1) \quad (66.67)$$

5.5.1. Calculating Probabilities

Property 66.16 Symmetry: Let $z > 0$

$$\mathbb{P}(Z \leq z) = \Phi(z) \quad (66.68)$$

$$\mathbb{P}(Z \leq -z) = \Phi(-z) = 1 - \Phi(z) \quad (66.69)$$

$$\mathbb{P}(-a \leq Z \leq b) = \Phi(b) - \Phi(-a) = \Phi(b) - (1 - \Phi(a))$$

$$a = b = z = 2\Phi(z) - 1 \quad (66.70)$$

5.5.2. Linear Transformations of Normal Dist.

$$\mathcal{N}(\mu, \sigma)$$

Proposition 66.1 [proof 66.12]

Linear Transformation:

Let X be a normally distributed random variable $X \sim \mathcal{N}(\mu, \sigma^2)$, then the linear transformed r.v. Y given by the affine transformation $Y = a + bX$ with $a \in \mathbb{R}, b \in \mathbb{R}_+$ follows:

$$Y \sim \mathcal{N}(a + b\mu, b^2\sigma^2) \iff f_Y(y) = \frac{1}{|b|} f_X\left(\frac{y-a}{b}\right) \quad (66.71)$$

Corollary 66.5

Linear Transformation from Standard Normal Dist.:

Let X be a standard normally distributed random variable $X \sim \mathcal{N}(0, 1)$, then the linear transformed r.v. Y given by the affine transformation $Y = a + bX$ with $a \in \mathbb{R}, b \in \mathbb{R}_+$ follows:

$$Y \sim \mathcal{N}(a, b^2) \iff f_Y(y) = \frac{1}{|b|} f_X\left(\frac{y-a}{b}\right) \quad (66.72)$$

Proposition 66.2 Standardization [proof 66.13]:

Let X be a normally distributed random variable $X \sim \mathcal{N}(\mu, \sigma^2)$, then there exists a linear transformation $Z = a + bX$ s.t. Z is a standard normally distributed random variable:

$$X \sim \mathcal{N}(\mu, \sigma^2) \xrightarrow{Z = \frac{X-\mu}{\sigma}} Z \sim \mathcal{N}(0, 1) \quad (66.73)$$

Note

If we know how many standard deviations our distribution is away from our target value then we can characterize it fully by the standard normal distribution.

Proposition 66.3 [proof 66.14]

Standardization of the CDF: Let $F_X(X)$ be the cumulative distribution function of a normally distributed random variable $X \sim \mathcal{N}(\mu, \sigma^2)$, then the cumulative distribution function $\Phi_Z(z)$ of the standardized random normal variable $Z \sim \mathcal{N}(0, 1)$ is related to $F_X(X)$ by:

$$F_X(x) = \Phi\left(\frac{x-\mu}{\sigma}\right) \quad (66.74)$$

6. The Multivariate Normal distribution

Definition 66.33

Multivariate Normaldistribution/Gaussian:

An \mathbb{R}^n -valued random variable $\mathbf{X} = (X_1, \dots, X_n)$ is **Multivariate Gaussian/Normaldistribution** if every linear combination of its components is a (one-dimensional) Gaussian:

$$\exists \mu, \Sigma : \mathcal{L}\left(\sum_{i=1}^n \alpha_i X_i\right) = \mathcal{N}(\mu, \Sigma) \quad \forall \alpha_i \in \mathbb{R} \quad (66.75)$$

(possible degenerated $\mathcal{N}(0, 0)$ for $\forall \alpha_j = 0$)

Note

- Joint vs. multivariate:** a joint normal distribution can be a multivariate normal distribution or a product of univariate normal distributions **but**
- Multivariate refers to the number of variables that are placed as inputs to a function.

Definition 66.34

Multivariate Normal distribution: $\mathbf{X} \sim \mathcal{N}_k(\mu, \Sigma)$

A k -dimensional random vector

$\mathbf{X} = (X_1, \dots, X_n)^\top$ with $\mu = (\mathbb{E}[x_1], \dots, \mathbb{E}[x_k])^\top$

and $k \times k$ p.s.d.covariance matrix:

$$\Sigma := \mathbb{E}[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^\top] = [\text{Cov}[x_i, x_j], 1 \leq i, j \leq k] \quad (66.76)$$

follows a k -dim multivariate normal/Gaussian distribution if its law^[def. 65.25] satisfies:

$$\mathbf{f}_{\mathbf{X}}(X_1, \dots, X_n) = \mathcal{N}(\mu, \Sigma) \quad (66.76)$$

$$= \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} \exp\left(-\frac{1}{2}(\mathbf{X} - \mu)^\top \Sigma^{-1} (\mathbf{X} - \mu)\right)$$

Normalisation

$$\varphi_{\mathbf{X}}(\mathbf{u}) = \exp\left\{i\mathbf{u}^\top \mu - \frac{1}{2}\mathbf{u}^\top \Sigma \mathbf{u}\right\} \quad (66.77)$$

Definition 66.35

$$\mathbf{X} \sim \mathcal{N}_k(\mu, \text{diag}(\sigma_1^2, \dots, \sigma_k^2))$$

Diagonal Gaussian Distribution

[proof 66.17]:

A diagonal Gaussian is a Multivariate Normaldistribution/Gaussian^[def. 66.34] with a diagonal covariance matrix with that can be decomposed into k independent distributions:

$$\mathbf{X} = (X_1, \dots, X_n)^\top \quad \text{with} \quad \mu = (\mathbb{E}[x_1], \dots, \mathbb{E}[x_k])^\top$$

and $k \times k$ p.s.d.covariance matrix:

$$\text{diag}(\sigma_1^2, \dots, \sigma_k^2)$$

and is given by:

$$\mathbf{f}_{\mathbf{X}}(X_1, \dots, X_k) = \mathcal{N}(\mu, \Sigma) = \prod_{i=1}^k f_{X_i}(X_i) \quad (66.78)$$

$$= \frac{1}{\sqrt{(2\pi)^k (\prod_{i=1}^k \sigma_i^2)}} \exp\left(-\frac{1}{2} \sum_{i=1}^k \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right)$$

Explanation 66.2 (Diagonal Gaussian Distribution). Is a Gaussian distribution that is scaled along the axis i.e. for a 2d distribution an ellipse along the x or y-axis.

Definition 66.36

$$\mathbf{X} \sim \mathcal{N}_k(\mu, \mathbf{I}_k \sigma^2)$$

Isotropic Gaussian

[proof 66.17]:

An isotropic Gaussian is a diagonal Multivariate Normaldistribution/Gaussian explanation 66.2 with constant standard deviation along the diagonal:

$$\mathbf{X} = (X_1, \dots, X_n)^\top \quad \text{with} \quad \mu = (\mathbb{E}[x_1], \dots, \mathbb{E}[x_k])^\top$$

and $k \times k$ p.s.d.covariance matrix:

$$\mathbf{I}_k \sigma = \text{diag}(1)_k \sigma = \begin{cases} \sigma & \text{if } i = j \\ 0 & \text{else} \end{cases}$$

and is given by:

$$\mathbf{f}_{\mathbf{X}}(X_1, \dots, X_k) = \mathcal{N}(\mu, \Sigma) = \prod_{i=1}^k f_{X_i}(X_i) \quad (66.79)$$

$$= \frac{1}{\sqrt{(2\pi\sigma^2)^k}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^k (x_i - \mu_i)^2\right)$$

1. Joint Gaussian Distributions

Definition 66.37 Jointly Gaussian Random Variables:

Two random variables X, Y both scalars or vectors, are said to be **jointly Gaussian** if the joint vector random variable $\mathbf{Z} = [X \ Y]^\top$ is again a GRV.

Property 66.17

[proof 66.16]

Joint Independent Gaussian Random Variables: Let X_1, \dots, X_n be \mathbb{R} -valued **independent** random variables with laws $\mathcal{N}(\mu_i, \sigma_i^2)$. Then the law of $\mathbf{X} = (X_1, \dots, X_n)^\top$ is a (multivariate) Gaussian distribution $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$ with:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix} \quad \text{and} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} \quad (66.80)$$

Corollary 66.6 Quadratic Form:

If \mathbf{x} and \mathbf{y} are both independent GRVs

$$\mathbf{x} \sim \mathcal{N}(\mu_x, \Sigma_x) \quad \mathbf{y} \sim \mathcal{N}(\mu_y, \Sigma_y)$$

then they are jointly Gaussian^[def. 66.37] given by:

$$\mathbf{p}(\mathbf{x}, \mathbf{y}) = \mathbf{p}(\mathbf{x})\mathbf{p}(\mathbf{y}) \quad (66.81)$$

$$\propto \exp\left(-\frac{1}{2}[(\mathbf{x} - \mu_x)^\top \Sigma_x^{-1}(\mathbf{x} - \mu_x) + (\mathbf{y} - \mu_y)^\top \Sigma_y^{-1}(\mathbf{y} - \mu_y)]\right)$$

$$= \exp\left(-\frac{1}{2}[(\mathbf{x} - \mu_x)^\top (\mathbf{y} - \mu_y)^\top \begin{bmatrix} \Sigma_x^{-1} & 0 \\ 0 & \Sigma_y^{-1} \end{bmatrix} [\mathbf{x} - \mu_x] \mathbf{y} - \mu_y]\right)$$

$$\cong \exp\left(-\frac{1}{2}(\mathbf{z} - \mu_z)^\top \Sigma_z^{-1}(\mathbf{z} - \mu_z)\right)$$

Property 66.18

Marginal Distribution of Multivariate Gaussian: Let $\mathbf{X} = (X_1 \dots X_n)^\top \sim \mathcal{N}(\mu, \Sigma)$ be a an \mathbb{R}^n valued Gaussian and let $V = \{1, 2, \dots, n\}$ be the index set of its variables. The k -variate marginal distribution of the Gaussian indexed by a subset of the variables:

$$A = \{i_1, \dots, i_k\} \quad i_j \in V \quad (66.82)$$

is given by:

$$\mathbf{X} = (X_{i_1} \dots X_{i_k})^\top \sim \mathcal{N}(\mu_A, \Sigma_{AA}) \quad (66.83)$$

$$\Sigma = \begin{bmatrix} \sigma_{i_1, i_1}^2 & \dots & \sigma_{i_1, i_k}^2 \\ \vdots & \ddots & \vdots \\ \sigma_{i_k, i_1}^2 & \dots & \sigma_{i_k, i_k}^2 \end{bmatrix} \quad \text{and} \quad \mu = \begin{bmatrix} \mu_{i_1} \\ \vdots \\ \mu_{i_k} \end{bmatrix}$$

2. Conditional Gaussian Distributions

Property 66.19 Conditional Gaussian Distribution: Let $\mathbf{X} = (X_1 \dots X_n)^\top \sim \mathcal{N}(\mu, \Sigma)$ be a an \mathbb{R}^n valued Gaussian and let $V = \{1, 2, \dots, n\}$ be the index set of its variables.

Suppose we take two disjoint subsets of V :

$$A = \{i_1, \dots, i_k\} \quad B = \{j_1, \dots, j_m\} \quad i_l, j_l \in V$$

then the conditional distribution of the random vector \mathbf{X}_A , conditioned on \mathbf{X}_B given by $\mathbb{P}(\mathbf{X}_A | \mathbf{X}_B = \mathbf{x}_B)$ is:

$$\mathbf{X}_A = (X_{i_1} \dots X_{i_k})^\top \sim \mathcal{N}(\mu_{A|B}, \Sigma_{A|B}) \quad (66.84)$$

$$\begin{aligned} \mu_{A|B} &= \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (\mathbf{x}_B - \mu_B) \\ \Sigma_{A|B} &= \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA} \end{aligned}$$

Note

Can be proofed using the matrix inversion lemma but is a very tedious computation.

Corollary 66.7

Conditional Distribution of Joint Gaussian's: Let \mathbf{X} and \mathbf{Y} be jointly Gaussian random vectors:

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix}\right) \quad (66.85)$$

then the marginal distribution of \mathbf{x} conditioned on \mathbf{y} can be written as:

$$\begin{aligned} \mathbf{x} &\sim \mathcal{N}(\mu_{X|Y}, \Sigma_{X|Y}) \\ \mu_{X|Y} &= \mu_X + \mathbf{C}\mathbf{B}^{-1}(\mathbf{y} - \mu_Y) \\ \Sigma_{X|Y} &= \mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top \end{aligned} \quad (66.86)$$

3. Transformations

Property 66.20 Multiples of Gaussian's: Let $\mathbf{X} = (X_1 \dots X_n)^\top \sim \mathcal{N}(\mu, \Sigma)$ be a an \mathbb{R}^n valued Gaussian and let $\mathbf{A} \in \mathbb{R}^{d \times n}$ then it follows:

$$Y = \mathbf{A}\mathbf{X} \in \mathbb{R} \quad Y \sim \mathcal{N}(\mathbf{A}\mu, \mathbf{A}\Sigma\mathbf{A}^\top) \quad (66.87)$$

Property 66.21 Affine Transformation of GRVs: Let $\mathbf{y} \in \mathbb{R}^n$ be GRV, $\mathbf{A} \in \mathbb{R}^{d \times n}$, $\mathbf{b} \in \mathbb{R}^d$ and let \mathbf{x} be defined by the affine transformation^[def. 59.45]:

$$\mathbf{x} = \mathbf{A}\mathbf{y} + \mathbf{b} \quad \mathbf{A} \in \mathbb{R}^{d \times n}, \mathbf{b} \in \mathbb{R}^d$$

Then \mathbf{x} is a GRV (see ?? 66.15).

Property 66.22 Linear Combination of jointly GRVs: Let $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$ two jointly GRVs, and let \mathbf{z} be defined as:

$$\mathbf{z} = \mathbf{A}_x \mathbf{x} + \mathbf{A}_y \mathbf{y} \quad \mathbf{A}_x \in \mathbb{R}^{d \times n}, \mathbf{A}_y \in \mathbb{R}^{d \times m}$$

Then \mathbf{z} is GRV (see ?? 66.18).

Definition 66.38 Gaussian Noise: Is statistical noise having a probability density function (PDF) equal to that of the normal/Gaussian distribution.

4. Gamma Distribution

$\Gamma(x, \alpha, \beta)$ Proofs

Definition 66.39 Gamma Distribution $X \sim \Gamma(x, \alpha, \beta)$: Is a widely used distribution that is related to the exponential distribution, Erlang distribution, and chi-squared distribution as well as Normal distribution:

$$f(x; \alpha, \beta) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (66.88)$$

$$\Gamma(\alpha) \stackrel{\text{eq. (54.81)}}{=} \int_0^\infty t^{\alpha-1} e^{-t} dt \quad (66.89)$$

with $\alpha, \beta \in \mathbb{R}_{>0}$

5. Chi-Square Distribution

6. Student's t-distribution

Definition 66.40 Student' t-distribution:

7. Delta Distribution

Definition 66.41 The delta function $\delta(\mathbf{x})$:

The delta/dirac function $\delta(\mathbf{x})$ is defined by:

$$\int_{\mathbb{R}} \delta(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = f(0)$$

for any integrable function f on \mathbb{R} .

Or alternatively by:

$$\delta(x - x_0) = \lim_{\sigma \rightarrow 0} \mathcal{N}(x|x_0, \sigma) \quad (66.90)$$

$$\approx \infty \mathbb{1}_{\{x=x_0\}} \quad (66.91)$$

Property 66.23 Properties of δ :

• **Normalization:** The delta function integrates to 1:

$$\int_{\mathbb{R}} \delta(x) dx = \int_{\mathbb{R}} \delta(x) \cdot c_1(x) dx = c_1(0) = 1$$

where $c_1(x) = 1$ is the constant function of value 1.

• **Shifting:**

$$\int_{\mathbb{R}} \delta(x - x_0) f(x) dx = f(x_0) \quad (66.92)$$

• **Symmetry:** $\int_{\mathbb{R}} \delta(-x) f(x) dx = f(0)$

• **Scaling:** $\int_{\mathbb{R}} \delta(a\alpha x) f(x) dx = \frac{1}{|\alpha|} f(0)$

Note

- In mathematical terms δ is not a function but a generalized function.
- We may regard $\delta(x - x_0)$ as a density with all its probability mass centered at the signle point x_0 .
- Using a box/indicator function s.t. its surface is one and its width goes to zero, instead of a normaldistribution eq. (66.90) would be a non-differentiable/discret form of the dirac measure.

Definition 66.42 Heaviside Step Function:

$$H(x) := \frac{d}{dx} \max\{x, 0\} \quad x \in \mathbb{R} \neq 0 \quad (66.93)$$

or alternatively:

$$H(x) := \int_{-\infty}^x \delta(s) ds \quad (66.94)$$

Proof 66.11 Definition 66.25: Consider a sequence of n random $\{X_i\}_{i=1}^n$ Bernoulli experiments^[def. 66.22] with success probability \mathbb{P} . Define the r.v. Y_n to be the sum of the n Bernoulli variables:

$$Y_n = \sum_{i=1}^n X_i \quad n \in \mathbb{N}$$

i.e. the total number of successes. Now lets calculate the probability density function f_n of Y_n . First let $(x_1, \dots, x_n) \in \{0, 1\}^n$ and let $y = \sum_{i=1}^n x_i$ a bit string of zeros and ones, with one occuring y times.

$$\begin{aligned} \mathbb{P}((X_1, X_2, \dots, X_n) = (x_1, x_2, \dots, x_n)) \\ = (\underbrace{\mathbb{P} \dots \mathbb{P}}_{y \text{-times}}) \cdot (\underbrace{\mathbb{P} \dots \mathbb{P}}_{n-y \text{-times}}) = \mathbb{P}^y (1 - \mathbb{P})^{n-y} \end{aligned}$$

However we need to take into account that there exists further realization $\mathbf{X} = \mathbf{x}$, that correspond to different orders of the elements in our two classes $\{0, 1\}$ which leads to $\frac{n!}{y!(n-y)!} = \binom{n}{y}$:

$$f_n(y) = \binom{n}{y} \mathbb{P}^y (1 - \mathbb{P})^{n-y} \quad y \in \{0, 1, \dots, n\}$$

Proof 66.12: proposition 66.1: Let X be normally distributed with $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$\begin{aligned} F_Y(y) \stackrel{y \geq 0}{=} \mathbb{P}_Y(Y \leq y) = \mathbb{P}(a + bX \leq y) = \mathbb{P}_X\left(X \leq \frac{y-a}{b}\right) \\ = \mathbb{P}_X\left(\frac{y-a}{b}\right) \end{aligned}$$

$$\begin{aligned} F_Y(y) \stackrel{y \leq 0}{=} \mathbb{P}_Y(Y \leq y) = \mathbb{P}(a + bX \leq y) = \mathbb{P}_X\left(X \geq \frac{y-a}{b}\right) \\ = 1 - \mathbb{P}_X\left(\frac{y-a}{b}\right) \end{aligned}$$

Differentiating both expressions w.r.t. y leads to:

$$f_Y(y) = \frac{1}{b} \frac{d\mathbb{P}_X\left(\frac{y-a}{b}\right)}{dy} = \frac{1}{b} \mathbb{P}'_X\left(\frac{y-a}{b}\right) = \frac{1}{b} f_X\left(\frac{y-a}{b}\right)$$

eq. (66.71)).

in order to prove that $Y \sim \mathcal{N}(a + b\mu, b^2\sigma^2)$ we simply plug f_X in the previous expression:

$$\begin{aligned} f_Y(y) &= \frac{1}{\sqrt{2\pi}\sigma|b|} \exp\left\{-\frac{1}{2} \left(\frac{y-a}{b} - \mu\right)^2\right\} \\ &= \frac{1}{\sqrt{2\pi}\sigma|b|} \exp\left\{-\frac{1}{2} \left(\frac{(y-a) - b\mu}{\sigma|b|}\right)^2\right\} \end{aligned}$$

Proof 66.13: proposition 66.2: Let X be normally distributed with $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$Z := \frac{X - \mu}{\sigma} = \frac{1}{\sigma} X - \frac{\mu}{\sigma} = aX + b \quad \text{with } a = \frac{1}{\sigma}, b = -\frac{\mu}{\sigma}$$

$$\text{eq. (66.71)} \quad \mathcal{N}(a\mu + b, a^2\sigma^2) \sim \mathcal{N}\left(\frac{\mu}{\sigma} - \frac{\mu}{\sigma}, \frac{\sigma^2}{\sigma^2}\right) \sim \mathcal{N}(0, 1)$$

Proof 66.14: proposition 66.3: Let X be normally distributed with $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$\begin{aligned} F_X(x) &= \mathbb{P}(X \leq x) \stackrel{-\mu}{\underset{\sigma}{\div}} \mathbb{P}\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) = \mathbb{P}\left(Z \leq \frac{x - \mu}{\sigma}\right) \\ &= \Phi\left(\frac{x - \mu}{\sigma}\right) \end{aligned}$$

Proof 66.15: Property 66.21 scalar case

Let $y \sim p(y) = \mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$ and define $\mathbf{x} = ay + b$ $a \in \mathbb{R}_+$, $b \in \mathbb{R}$. Using the Change of variables formula it follows:

$$\begin{aligned} p_x(\bar{x}) &\stackrel{\text{eq. (65.46)}}{=} \frac{p(y)}{|\frac{dy}{dx}|} = \frac{\frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(\bar{x}-a)^2}{2\sigma^2}\right)}{|\frac{d\bar{x}}{dy}|} = \frac{\frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(\bar{x}-a)^2}{2\sigma^2}\right)}{\left|\frac{a}{\sigma}\right|} \\ &= \frac{1}{\sqrt{2\pi}a^2\sigma^2} \exp\left(-\frac{(\bar{x}-a)^2}{2a^2\sigma^2}\right) \end{aligned}$$

Hence $x \sim \mathcal{N}(\mu_x, \sigma_x^2) = \mathcal{N}(a\mu + b, a^2\sigma^2)$

Note

We can also verify that we have calculated the right mean and variance by:

$$\mathbb{E}[x] = \mathbb{E}[ay + b] = a\mathbb{E}[y] + b = a\mu + b$$

$$\mathbb{V}[x] = \mathbb{V}[ay + b] = a^2\mathbb{V}[y] = a^2\sigma^2$$

Proof 66.16:

$$\begin{aligned} p_{\mathbf{X}}(\mathbf{u}) &= \prod_i^n p_{X_i}(u_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i^2} \exp\left(-\frac{(u_i - \mu_i)^2}{2\sigma_i^2}\right) \\ \varphi_{\mathbf{X}}(\mathbf{u}) &= \exp\left\{iu_1\mu_1 - \frac{1}{2}\mathbf{u}_1\mathbf{u}_1^\top\right\} \dots \exp\left\{iu_n\mu_n - \frac{1}{2}\mathbf{u}_n\mathbf{u}_n^\top\right\} \\ &= \exp\left\{i \sum_i^n u_i \mu_i - \frac{1}{2} \sum_i^n \sigma_i u_i^2\right\} = \exp\left\{i\mathbf{u}^\top \mu - \frac{1}{2}\mathbf{u}^\top \Sigma \mathbf{u}\right\} \end{aligned}$$

Proof 66.17 Diagonal Gaussian Distribution^[def. 66.35]:

$$\begin{aligned} \Sigma^{-1} &= \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sigma_k^2} \end{bmatrix} \quad |\Sigma^{-1}| = \prod_{i=1}^k \sigma_i^2 = \left(\prod_{i=1}^k \sigma_i\right)^2 \\ (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) &= \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sigma_k^2} \end{bmatrix} \begin{bmatrix} (x_1 - \mu_1) \\ (x_2 - \mu_2) \\ \vdots \\ (x_k - \mu_k) \end{bmatrix} \end{aligned}$$

$$= [(x_1 - \mu_1) \ (x_2 - \mu_2) \ \dots \ (x_k - \mu_k)] \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sigma_k^2} \end{bmatrix} \begin{bmatrix} (x_1 - \mu_1) \\ (x_2 - \mu_2) \\ \vdots \\ (x_k - \mu_k) \end{bmatrix}$$

$$= \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} + \dots + \frac{(x_k - \mu_k)^2}{\sigma_k^2} = \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}$$

Combining those two lead directly to:

$$f_{\mathbf{X}}(x_1, \dots, x_k) = \frac{1}{\sqrt{(2\pi)^k (\prod_{i=1}^k \sigma_i^2)}} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right)$$

Proof 66.18: Property 66.22

From Property 66.21 it follows immediately that \mathbf{z} is GRV

$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_z, \Sigma_z)$ with:

$$\mathbf{z} = \mathbf{A}\xi \quad \text{with} \quad \mathbf{A} = [\mathbf{A}_x \quad \mathbf{A}_y] \quad \text{and} \quad \xi = (\mathbf{x} \quad \mathbf{y})$$

Knowing that \mathbf{z} is a GRV it is sufficient to calculate $\boldsymbol{\mu}_z$ and Σ_z in order to characterize its distribution:

$$\mathbb{E}[\mathbf{z}] = \mathbb{E}[\mathbf{A}_x \mathbf{x} + \mathbf{A}_y \mathbf{y}] = \mathbf{A}_x \boldsymbol{\mu}_x + \mathbf{A}_y \boldsymbol{\mu}_y$$

$$\begin{aligned} \mathbb{V}[\mathbf{z}] &= \mathbb{V}[\mathbf{A}\xi] = \mathbf{A}\mathbb{V}[\xi]\mathbf{A}^\top \quad ?? \\ &= [\mathbf{A}_x \quad \mathbf{A}_y] \begin{bmatrix} \mathbb{V}[x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \mathbb{V}[y] \end{bmatrix} [\mathbf{A}_x \quad \mathbf{A}_y]^\top \\ &= [\mathbf{A}_x \quad \mathbf{A}_y] \begin{bmatrix} \mathbb{V}[x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \mathbb{V}[y] \end{bmatrix} \begin{bmatrix} \mathbf{A}_x^\top \\ \mathbf{A}_y^\top \end{bmatrix} \\ &= \mathbf{A}_x \mathbb{V}[x] \mathbf{A}_x^\top + \mathbf{A}_y \mathbb{V}[y] \mathbf{A}_y^\top \\ &\quad + \underbrace{\mathbf{A}_y \text{Cov}[y, x] \mathbf{A}_x^\top}_{=0 \text{ by independence}} + \underbrace{\mathbf{A}_x \text{Cov}[x, y] \mathbf{A}_y^\top}_{=0 \text{ by independence}} \\ &= \mathbf{A}_x \Sigma_x \mathbf{A}_x^\top + \mathbf{A}_y \Sigma_y \mathbf{A}_y^\top \end{aligned}$$

Note

Can also be proofed by using the normal definition of [def. 66.15] and tedious computations.

Proof 66.19: Equation (66.43) If $\mathbf{x} = c_i$ i.e. the outcome c_i has occurred then it follows:

$$\prod_j^k p_i^{\delta[x=c_i]} = p_1^0 \cdots p_i^1 \cdots p_k^0 = 1 \cdots p_i \cdots 1 = p(\mathbf{x} = c_i | \mathbf{p})$$

Sampling Methods

1. Sampling Random Numbers

Most math libraries have uniform **random number generator (RNG)** i.e. functions to generate uniformly distributed random numbers $U \sim \mathcal{U}[a, b]$ (eq. (66.52)).

Furthermore repeated calls to these RNG are independent, that is:

$$\begin{aligned} p_{U_1, U_2}(u_1, u_2) &\stackrel{??}{=} p_{U_1}(u_1) \cdot p_{U_2}(u_2) \\ &= \begin{cases} 1 & \text{if } u_1, u_2 \in [a, b] \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Question: using samples $\{u_1, \dots, u_n\}$ of these CRVs with uniform distribution, how can we create random numbers with arbitrary discrete or continuous PDFs?

2. Inverse-transform Technique

Idea

Can make use of section 1 and the fact that CDF are increasing functions (^[def. 54.12]). **Advantage:**

- Simple to implement
- All discrete distributions can be generated via inverse-transform technique

Drawback:

- Not all continuous distributions can be integrated/have closed form solution for their CDF.
- E.g. Normal-, Gamma-, Beta-distribution.

1. Continuous Case

Definition 67.1 One Continuous Variable: Given: a desired continuous pdf f_X and uniformly distributed rn $\{u_1, u_2, \dots\}$:

1. Integrate the desired pdf f_X in order to obtain the desired cdf F_X :

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \quad (67.1)$$

2. Set $F_X(X) = U$ on the range of X with $U \sim \mathcal{U}[0, 1]$.

3. Invert this equation/find the inverse $F_X^{-1}(U)$ i.e. solve:

$$U = F_X(X) = F_X\left(\underbrace{\underline{F}_X^{-1}(U)}_{X}\right) \quad (67.2)$$

4. Plug in the uniformly distributed rn:

$$x_i = \underline{F}_X^{-1}(u_i) \quad \text{s.t.} \quad x_i \sim f_X \quad (67.3)$$

Definition 67.2 Multiple Continuous Variables:

Given: a pdf of multiple rvs $f_{X,Y}$:

1. Use the product rule $(??)$ in order to decompose $f_{X,Y}$:

$$f_{X,Y} = f_{X|Y}(x, y) = f_{X|Y}(x|y)f_Y(y) \quad (67.4)$$

2. Use ^[def. 67.3] to first get a rv for y of $Y \sim f_Y(y)$.

3. Then with this fixed y use ^[def. 67.3] again to get a value for x of $X \sim f_{X|Y}(x|y)$.

Proof 67.1: ^[def. 67.3]:

Claim: if U is a uniform rv on $[0, 1]$ then $\underline{F}_X^{-1}(U)$ has F_X as its CDF.

Assume that F_X is strictly increasing (^[def. 54.12]).

Then for any $u \in [0, 1]$ there must exist a unique x s.t. $F_X(x) = u$.

Thus F_X must be invertible and we may write $x = \underline{F}_X^{-1}(u)$.

Now let a arbitrary:

$$F_X(a) = \mathbb{P}(\underline{x} \leq a) = \mathbb{P}(\underline{F}_X^{-1}(U) \leq a)$$

Since F_X is strictly increasing:

$$\begin{aligned} \mathbb{P}(\underline{F}_X^{-1}(U) \leq a) &= \mathbb{P}(U \leq F_X(a)) \\ &\stackrel{\text{eq. (66.52)}}{=} \int_0^{F_X(a)} 1 dt = F_X(a) \end{aligned}$$

Note

Strictly speaking we may not assume that a CDF is **strictly** increasing but we as all CDFs are weakly increasing (^[def. 54.12]) we may always define an auxiliary function by its infimum:

$$\hat{F}_X^{-1} := \inf \{x | F_X(x) \geq 0\} \quad u \in [0, 1] \quad (67.5)$$

2. Discret Case

Idea

Given: a desired $U \sim \mathcal{U}[0, 1]$ discret pmf p_X s.t. 1
 $\mathbb{P}(X = x_i) = p_X(x_i)$ and uniformly distributed rn $\{u_1, u_2, \dots\}$.
Goal: given a uniformly distributed rn u determine k s.t.:

$$\sum_{i=1}^{k-1} < U \leq \sum_{i=1}^k \iff F_X(x_{k-1}) < u \leq F_X(x_k) \quad (67.6)$$

and return x_k .

Definition 67.3 One Discret Variable:

1. Compute the CDF of p_X (^[def. 66.8])

$$F_X(x) = \sum_{t=-\infty}^x p_X(t) \quad (67.7)$$

2. Given the uniformly distributed rn $\{u_i\}_{i=1}^n$ find k^i (\cong inversion) s.t.:

$$F_X(x_{k(i)-1}) < u_i \leq F_X(x_{k(i)}) \quad \forall u_i \quad (67.8)$$

Proof 67.2: ???: First of all notice that we can always solve for an unique x_k . Given a fixed x_k determine the values of u for which:

$$F_X(x_{k-1}) < u \leq F_X(x_k) \quad (67.9)$$

Now observe that:

$$\begin{aligned} u &\leq F_X(x_k) = F_X(x_{k-1}) + p_X(x_k) \\ &\Rightarrow F_X(x_{k-1}) < u \leq F_X(x_{k-1}) + p_X(x_k) \end{aligned}$$

The probability of U being in $(F_X(x_{k-1}), F_X(x_k)]$ is:

$$\begin{aligned} \mathbb{P}(U \in [F_X(x_{k-1}), F_X(x_k)]) &= \int_{F_X(x_{k-1})}^{F_X(x_k)} p_U(t) dt \\ &= \int_{F_X(x_{k-1})}^{F_X(x_k)} 1 dt = \int_{F_X(x_{k-1})}^{F_X(x_{k-1})+p_X(x_k)} 1 dt = p_X(x_k) \end{aligned}$$

Hence the random variable $x_k \in \mathcal{X}$ has the pdf p_X .

Definition 67.4

Multiple Continuous Variables (Option 1):

Given: a pdf of multiple rvs $p_{X,Y}$:

1. Use the product rule $(??)$ in order to decompose $p_{X,Y}$:

$$p_{X,Y} = p_{X|Y}(x, y) = p_{X|Y}(x|y)p_Y(y) \quad (67.10)$$

2. Use $??$ to first get a rv for y of $Y \sim f_Y(y)$.

3. Then with this fixed y use $??$ again to get a value for x of $X \sim p_{X|Y}(x|y)$.

Definition 67.5

Multiple Continuous Variables (Option 2):

Note: this only works if \mathcal{X} and \mathcal{Y} are finite.

Given: a pdf of multiple rvs $p_{X,Y}$ let $N_x = |\mathcal{X}|$ and $N_y = |\mathcal{Y}|$ the number of elements in \mathcal{X} and \mathcal{Y} .

Define $p_Z(1) = p_{X,Y}(1, 1), p_Z(2) = p_{X,Y}(1, 2), \dots, p_Z(N_x \cdot N_y) = p_{X,Y}(N_x, N_y)$

Then simply apply $??$ to the auxillary pdf p_Z

1. Use the product rule $(??)$ in order to decompose $p_{X,Y}$:

$$f_{X,Y} = f_{X|Y}(x, y) = f_{X|Y}(x|y)f_Y(y) \quad (67.11)$$

2. Use ^[def. 67.3] to first get a rv for y of $Y \sim f_Y(y)$.

3. Then with this fixed y use ^[def. 67.3] again to get a value for x of $X \sim f_{X|Y}(x|y)$.

3. Monte Carlo Methods

1. Monte Carlo (MC) Integration

Integration methods s.a. Simpson integration ^[def. 62.34] suffer heavily from the curse of dimensionality.

An n-order ^[def. 62.31] quadrature scheme \mathcal{Q}_n in 1-dimension is usually of order n/d in d-dimensions.

Idea estimate an integral stochastically by drawing sample from some distribution.

Definition 67.6 Monte Carlo Integration:

$$3 + 4 \quad (67.12)$$

2. Rejection Sampling

3. Importance Sampling

Descriptive Statistics

1. Populations and Distributions

Definition 68.1 Population $\{x_i\}_{i=1}^N$: is the entire set of entities from which we can draw sample.

Definition 68.2 Families of Probability Distributions p_{θ} : Are probability distributions that vary only by a set of hyper parameters θ [def. 68.1].

Definition 68.3 Population/Statistical Parameter θ : Are the parameters defining families of probability distributions [def. 68.2].

Explanation 68.1 (Definition 68.1). Such hyper parameters are often characterized by populations following a certain family of distributions with the help of a statistic. Hence they are called population or statistical parameters.

1. Characteristics of Populations

Definition 68.4 Population Mean: Given a population $\{x_i\}_{i=1}^N$ of size N its variance is defined as:

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad (68.1)$$

Definition 68.5 Population Variance: Given a population $\{x_i\}_{i=1}^N$ of size N its variance is defined as: $\{\hat{x}_i\}_{i=1}^N$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \quad (68.2)$$

Note

The population variance and mean are equally to the mean derived from the true distribution of the population.

2. Sample Statistics

Definition 68.6 (Sample) Statistic: A statistic is a measurable function T that assigns a single value t to a sample of random variables or population:

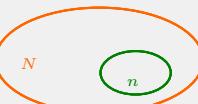
$$t : \mathbb{R}^n \mapsto \mathbb{R} \quad t = T(X_1, \dots, X_n)$$

E.g. T could be the mean, variance,...

Definition 68.7 Degrees of freedom of a Statistic: Is the number of values in the final calculation of a statistic that are free to vary.

Note

The function itself is independent of the sample's distribution; that is, the function can be stated before realization of the data.



3. Point and Interval Estimation

Assume a population X with a given sample $\{x_i\}_{i=1}^n$ follows some family of distributions:

$$X \sim p_X(\cdot; \theta) \quad (68.3)$$

how can we estimate the correct value of the parameter θ or some function of that parameter $\tau(\theta)$?

1. Point Estimates

Definition 68.8 (Point) Estimator $\hat{\theta}$: Is a statistic [def. 68.6] that tries estimates an unknown parameter θ of an underlying family of distributions [def. 68.2] for a given sample $\{x_i\}_{i=1}^n$ of that distribution:

$$\hat{\theta} = t(x_1, \dots, x_n) \quad (68.4)$$

Note

The other kind of estimators are interval estimators which do not calculate a statistic but an interval of plausible values of an unknown population parameter θ .

The most prevalent forms of interval estimation are:

- Confidence intervals (frequentist method).
- Credible intervals (Bayesian method).

3.1.1. Empirical Mean

Definition 68.9 Sample/Empirical Mean \bar{x} : The sample mean is an estimate/statistic of the population mean [def. 68.4] and can be calculated from an observation/sample of the total population $\{x_i\}_{i=1}^n \subset \{x_i\}_{i=1}^N$:

$$\bar{x} = \hat{\mu}_X = \frac{1}{n} \sum_{i=1}^n x_i \quad (68.5)$$

Corollary 68.1 [proof 68.1]

Unbiased Sample Mean:

The sample mean estimator is unbiased:

$$\mathbb{E}[\hat{\mu}_X] = \mu \quad (68.6)$$

Corollary 68.2 [Proof 68.2]

Variance of the Sample Mean:

The variance of the sample mean estimator is given by:

$$\text{Var}[\hat{\mu}_X] = \frac{1}{n} \sigma_X^2 \quad (68.7)$$

3.1.2. Empirical Variance

Definition 68.10 Biased Sample Variance:

The sample variance is an estimate/statistic of the population variance [def. 68.5] and can be calculated from an observation/sample of the total population $\{x_i\}_{i=1}^n \subset \{x_i\}_{i=1}^N$:

$$s^2 = \hat{\sigma}_X^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (68.8)$$

Definition 68.11 [proof 68.3]

(Unbiased) Sample Variance:

The unbiased form of the sample variance [def. 68.10] is given by:

$$s^2 = \hat{\sigma}_X^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (68.9)$$

Definition 68.12 Bessel's Correction: The factor

$$\frac{n}{n-1} \quad (68.10)$$

is called Bessel's correction. Multiplying the uncorrected population variance eq. (68.8) by this term yields an unbiased estimated of the variance.

Attention:

- The Bessel correction holds for the variance but not for the standard deviation.
- Usually only the unbiased variance is used and sometimes also denoted by s_n^2

2. Interval Estimates

Definition 68.13 Interval Estimator

$\hat{\theta}$: Is an estimator that tries to bound an unknown parameter θ of an underlying family of distributions [def. 68.2] for a given sample $\{x_i\}_{i=1}^n$ of that distribution.

Let $\theta \in \Theta$ and define two point statistics [def. 68.6] g and h then an interval estimate is defined as:

$$\mathbb{P}(L_n < \theta < U_n) = \gamma \quad \forall \theta \in \Theta \quad L_n = g(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ \gamma \in [0, 1] \quad U_n = h(\mathbf{x}_1, \dots, \mathbf{x}_n) \quad (68.11)$$

Statistical Tests

4. Parametric Hypothesis Testing

Definition 68.14 Parametric Hypothesis Testing:

Hypothesis testing is a statistical procedure in which a hypothesis is tested based on sampled data X_1, \dots, X_n .

1. Null Hypothesis

Definition 68.15 Null Hypothesis H_0 : A null hypothesis H_0 is an assumption on a population [def. 68.1] parameter [def. 68.3] θ :

$$H_0 : \theta = \theta_0 \quad (68.12)$$

Note

Often, a null hypothesis cannot be verified, but can only be falsified.

Definition 68.16 Alternative Hypothesis H_A/H_1 : The alternative hypothesis H_1 is an assumption on a population [def. 68.1] parameter [def. 68.3] θ that is opposite to the null hypothesis.

$$H_A : \theta \begin{cases} > \theta_0 & (\text{one-sided}) \\ < \theta_0 & (\text{one-sided}) \\ \neq \theta_0 & (\text{two-sided}) \end{cases} \quad (68.13)$$

2. Test Statistic

The decision on the hypothesis test is based on a sample from the population $X(n) = \{X_1, \dots, X_n\}$ however the decision is usually not based on single sample but a sample statistic [def. 68.6] as this is easier to use.

Definition 68.17 [example 68.4]

Test Statistic/Testing Parameter

T : Is a sample statistic [def. 68.6] used for hypothesis tests in order to give evidence for or against a hypothesis:

$$t_n = T(D_n) = T(\{X_1, \dots, X_n\}) \quad (68.14)$$

3. Sampling Distribution

Definition 68.18 $T_{\theta_0}(t)$

Null Distribution/Sampling Distribution under H_0 : Let $D_n = \{X_1, \dots, X_n\}$ be a random sample from the true population p_{pop} and let $T(D_n)$ be a test statistic of that sample.

The probability distribution of the test statistic under the assumption that the null hypothesis is true is called *sampling distribution*:

$$t \sim T_{\theta_0}(t) \quad X_i \sim p_{\text{pop}} \quad (68.15)$$

4. The Critical Region

Given a sample $D_n = \{X_1, \dots, X_n\}$ of the true population p_{pop} how should we decide whether the null hypothesis should be rejected or not?

Idea: let \mathcal{T} be the set of all possible values that the sample statistic T can map to. Now lets split \mathcal{T} in two disjoint sets \mathcal{T}_0 and \mathcal{T}_1 :

$$\mathcal{T} = \mathcal{T}_0 \cup \mathcal{T}_1 \quad \mathcal{T}_0 \cap \mathcal{T}_1 = \emptyset$$

- if $t_n = T(X_n) \in \mathcal{T}_0$ we accept the null hypothesis H_0
- if $t_n = T(X_n) \in \mathcal{T}_1$ we reject the null hypothesis for H_1

Definition 68.19 Critical/Rejection Region \mathcal{T}_1 : Is the set of all values of the test statistic [def. 68.17] t_n that causes us to reject the Null Hypothesis in favor for the alternative hypothesis H_A :

$$K = \mathcal{T}_1 = \{\mathcal{T} : H_0 \text{ rejected}\} \quad (68.16)$$

Definition 68.20 Acceptance Region \mathcal{T}_0 : Is the region where we accept the null hypothesis H_0 .

$$\mathcal{T}_0 = \{\mathcal{T} : H_0 \text{ accepted}\} \quad (68.17)$$

Definition 68.21 Critical Value c : Is the value of the critical region $c \in \mathcal{T}_1$ which is closest to the region of acceptance [def. 68.20]:

5. Type I&II Errors

Definition 68.22

False Positive

Type I Error: Is the rejection of the null hypothesis H_0 , even-tough it is true

$$\text{Test rejects } H_0 | H_0 \text{ true} \\ \iff t_n \in \mathcal{T}_1 | H_0 \text{ true} \quad (68.18)$$

Definition 68.23 False Negative

Type II Error:

Is the acceptance of a null hypothesis H_0 , even-tough its false:
Test accepts $H_0 | H_A$ true
 $\iff t_n \in \mathcal{T}_0 | H_A$ true

Types of Errors

Decision	H_0 true	H_0 false
Accept	TN	Type II (FN)
Reject	Type I (FP)	TP

6. Statistical Significance & Power

Question: how should we choose the split $\{\mathcal{T}_0, \mathcal{T}_1\}$? The bigger we choose \mathcal{T}_1 (and thus the smaller \mathcal{T}_0) the more likely it is to accept the alternative.

Idea: take the position of the adversary and choose \mathcal{T}_1 so small that $\theta \in \mathcal{T}_1$ has only a small probability of occurring.

Definition 68.24 [example 68.5] (Statistical) Significance α : A study's defined significance level α denotes the probability to incur a Type I Error [def. 68.22]:

$$\mathbb{P}(t_n \in \mathcal{T}_1 | H_0 \text{ true}) = \mathbb{P}(\text{test rejects } H_0 | H_0 \text{ true}) \leq \alpha \quad (68.20)$$

Definition 68.25 Probability Type II Error β : A test probability to for a false negative [def. 68.23] is defined as:

$$\mathbb{P}(t_n \in \mathcal{T}_0 | H_1 \text{ true}) = \mathbb{P}(\text{test accepts } H_0 | H_1 \text{ true}) \quad (68.21)$$

Definition 68.26 (Statistical) Power

1 - β :

A study's power $1 - \beta$ denotes a tests probability for a true positive:

$$1 - \beta(t_n) = \mathbb{P}(t_n \in \mathcal{T}_1 | H_1 \text{ true}) = \mathbb{P}(\text{test rejects } H_0 | H_1 \text{ true}) \quad (68.22)$$

$$1 - \beta(t_n) = \mathbb{P}(t_n \in \mathcal{T}_0 | H_1 \text{ true}) = \mathbb{P}(\text{test accepts } H_0 | H_1 \text{ true}) \quad (68.23)$$

Corollary 68.3 Types of Split: The Critical region is chosen s.t. we incur a Type I Error with probability less than α , which corresponds to the type of the test [def. 68.16]:

$$\begin{aligned} \mathbb{P}(c_2 \leq X \leq c_1) &\leq \alpha && \text{two-sided} \\ \text{or} \quad \mathbb{P}(c_2 \leq X) &\leq \frac{\alpha}{2} && \text{and} \quad \mathbb{P}(X \leq c_1) \leq \frac{\alpha}{2} \\ \mathbb{P}(c_2 \leq X) &\leq \alpha && \text{one-sided} \\ \mathbb{P}(X \leq c_1) &\leq \alpha && \text{one-sided} \end{aligned}$$

Truth	H_0 true	H_0 false
Decision		
H_0 accept	$1 - \alpha$	$1 - \beta$
H_0 rejected	α	β

7. P-Value

Definition 68.27 P-Value p : Given a test statistic $t_n = T(X_1, \dots, X_n)$ the p-value $p \in [0, 1]$ is the smallest value s.t. we reject the null hypothesis:

$$p := \inf \{\alpha | t_n \in \mathcal{T}_1\} \quad t_n = T(X_1, \dots, X_n) \quad (68.24)$$

Explanation 68.2.

The smaller the p-value the less likely is an observed statistic t_n and thus the higher is the evidence against a null hypothesis.

A null hypothesis has to be rejected if the p-value is bigger than the chosen significance niveau α .

5. Conducting Hypothesis Tests

- ① Select an appropriate test statistic^[def. 68.17] T .
- ② Define the null hypothesis H_0 and the alternative hypothesis H_A for T .
- ③ Find the sampling distribution^[def. 68.18] $T_{\theta_0}(t)$ for T , given H_0 true.
- ④ Choose the significance level α .
- ⑤ Evaluate the test statistic $t_n = T(X_1, \dots, X_n)$ for the sampled data.
- ⑥ Determine the p-value p .
- ⑦ Make a decision (accept or reject H_0)

1. Tests for Normally Distributed Data

Let us consider an i.i.d. sample of observations $\{x_i\}_{i=1}^n$, of a normally distributed population $X_{\text{pop}} \sim \mathcal{N}(\mu, \sigma^2)$.

From eqs. (68.6) and (68.7) it follows that the *mean of the sample* is distributed as:

$$\bar{X}_n \sim \mathcal{N}(\mu, \sigma^2/n)$$

thus the mean of the sample \bar{X}_n should equal the mean μ of the population. We now want to test the null hypothesis:

$$H_0 : \mu = \mu_0 \iff \bar{X}_n \sim \mathcal{N}(\mu_0, \sigma^2/n) \quad (68.25)$$

This is obviously only likely if the realization \bar{x}_n is close to μ_0 .

5.1.1. Z-Test σ known

Definition 68.28 Z-Test:

For a realization of Z with $\{x_i\}_{i=1}^n$ and mean \bar{x}_n :

$$z = \frac{\bar{x}_n - \mu_0}{\sigma/\sqrt{n}}$$

we reject the null hypothesis $H_0 : \mu = \mu_0$ for the alternative H_A for significance niveau^[def. 68.24] α if:

$$\begin{aligned} |z| \geq z_{1-\frac{\alpha}{2}} &\iff z \leq z_{\frac{\alpha}{2}} \vee z \geq z_{1-\frac{\alpha}{2}} \\ &\iff z \in \mathcal{T}_1 = \left(-\infty, -z_{\frac{\alpha}{2}}\right] \cup \left[z_{1-\frac{\alpha}{2}}, \infty\right) \\ z \geq z_{1-\alpha} &\iff z \in \mathcal{T}_1 = [z_{1-\alpha}, \infty) \\ z \leq z_{\alpha} = -z_{1-\alpha} &\iff z \in \mathcal{T}_1 = (-\infty, -z_{\alpha}] = (\infty, -z_{1-\alpha}] \end{aligned} \quad (68.26)$$

Notes

- Recall from [def. 66.19] and [cor. 66.4] that:
i.e. $\alpha = 0.05 \iff z_{0.05} = \Phi^{-1}(\alpha) \iff P(Z \leq z_{0.05}) = 0.05$
- $|z| \geq z_{1-\frac{\alpha}{2}}$ which stands for:

$$P(Z \leq z_{0.05}) + P(Z \geq z_{0.95}) = P(Z \leq -z_{1-0.05}) + P(Z \geq z_{0.95}) = P(|Z| \geq z_{0.95})$$

can be rewritten as:

$$z \geq z_{1-\frac{\alpha}{2}} \vee -z \geq z_{1-\frac{\alpha}{2}} \iff z \leq -z_{1-\frac{\alpha}{2}} = z_{\frac{\alpha}{2}}$$

- One usually goes over to the standard normal distribution proposition 66.2 and thus test how far one is away from zero mean \Rightarrow Z-test.
- We thus inquire a Type I error with probability α and should be small i.e. 1%.

5.1.2. t-Test σ unknown

In reality we usually do not know the true σ of the whole data set and thus calculate it over our sample. This however increases uncertainty and thus our sample does no longer follow a normal distribution but a **t-distribution** with $n-1$ degrees of freedom:

$$T \sim t_{n-1} \quad (68.27)$$

Definition 68.29 t-Test:

For a realization of T with $\{x_i\}_{i=1}^n$ and mean \bar{x}_n :

$$t = \frac{\bar{x}_n - \mu_0}{s_n/\sqrt{n}}$$

we reject the null hypothesis $H_0 : \mu = \mu_0$ for the alternative H_A if:

$$\begin{aligned} |t| &\geq t_{n-1, 1-\frac{\alpha}{2}} \\ &\iff t \in \mathcal{T}_1 = \left(-\infty, -t_{n-1, 1-\frac{\alpha}{2}}\right] \cup \left[t_{n-1, 1-\frac{\alpha}{2}}, \infty\right) \\ t &\geq t_{n-1, 1-\alpha} \\ &\iff t \in \mathcal{T}_1 = [t_{n-1, 1-\alpha}, \infty) \\ t &\leq t_{n-1, \alpha} = -t_{n-1, 1-\alpha} \\ &\iff t \in \mathcal{T}_1 = (-\infty, -t_{n-1, \alpha}] = (\infty, -t_{n-1, 1-\alpha}] \end{aligned}$$

Notes

- The t-distribution has fatter tails as the normal distribution \Rightarrow rare event become more likely
- For $n \rightarrow \infty$ the t-distribution goes over into the normal distribution
- The t-distribution gains a degree of freedom for each sample and loses one for each parameter we are interested in \Rightarrow n -samples and we are interested in one parameter μ .

2. Confidence Intervals

Now we are interested in the opposite of the critical region^[def. 68.19] namely the region of plausible values.

Definition 68.30 Confidence Interval

I:

Let $D_n = \{X_1, \dots, X_n\}$ be a sample of observations and T_n a sample statistic of that sample. The confidence interval is defined as:

$$I(D_n) = \{\theta_0 : T_n(D_n) \in \mathcal{T}_0\} = \{\theta_0 : H_0 \text{ is not rejected}\} \quad (68.28)$$

Corollary 68.4 : The confidence interval captures the unknown parameter θ with probability $1 - \alpha$:

$$P_\theta(\theta \in I(D_n)) = P(T_n(D_n) \in \mathcal{T}_0) = 1 - \alpha \quad (68.29)$$

6. Inferential Statistics

Goal of Inference

- ① What is a good guess of the parameters of my model?
- ② How do I quantify my uncertainty in the guess?

7. Examples

Example 68.1 ??: Let x be uniformly distributed on $[0, 1]$ (^{def. 66.28}) with pmf $p_X(x)$ then it follows:

$$\frac{dy}{dx} = \frac{1}{p_Y(y)} \Rightarrow dx = dy p_Y(y) \Rightarrow x = \int_{-\infty}^y p_Y(t) dt = F_Y(x)$$

Example 68.2 ??: Let

Example 68.3 Family of Distributions: The family of normal distribution \mathcal{N} has two parameters $\{\mu, \sigma^2\}$

Example 68.4 Test Statistic: Lets assume the test statistic follows a normal distribution:

$$T \sim \mathcal{N}(\mu; 1)$$

however we are unsure about the population parameter ^{def. 68.3} $\theta = \mu$ but assume its equal to θ_0 thus the null-and alternative hypothesis are:

$$H_0 : \mu = \mu_0 \quad H_1 : \mu \neq \mu_0$$

Example 68.5 Binomialtest:

Given: a manufacturer claims that a maximum of 10% of its delivered components are substandard goods.

In a sample of size $n = 20$ we find $x = 5$ goods that do not fulfill the standard and are skeptical that what the manufacturer claims is true, so we want to test:

$$H_0 : p = p_0 = 0.1 \quad \text{vs.} \quad H_A : p > 0.1$$

We model the number of number of defective goods using the binomial distribution ^{def. 66.25}

$$X \sim \mathcal{B}(n, p), n = 20 \quad P(X \geq x) = \sum_{k=x}^n \binom{n}{k} p^k (1-p)^{n-k}$$

from this we find:

$$\begin{aligned} P_{p_0}(X \geq 4) &= 1 - P_{p_0}(X \leq 3) = 0.13 \\ P_{p_0}(X \geq 5) &= 1 - P_{p_0}(X \leq 4) = 0.04 \leq \alpha \end{aligned}$$

thus the probability that equal 5 or more then 5 parts out of the 20 are rejects is less then 4%.

⇒ throw away null hypothesis for the 5% niveau in favor to the alternative.

⇒ the 5% significance niveau is given by $K = \{5, 6, \dots, 20\}$

Note

If $x < n/2$ it is faster to calculate $P(X \geq x) = 1 - P(X \leq x-1)$

8. Proofs

Proof 68.1: ^[cor. 68.1]

$$\mathbb{E}[\hat{\mu}_X] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n x_i\right] = \frac{1}{n} \mathbb{E}[\underbrace{\mu + \dots + \mu}_{1, \dots, n}]$$

Proof 68.2: ^[cor. 68.2]

$$\begin{aligned} \mathbb{V}[\hat{\mu}_X] &= \mathbb{V}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] \stackrel{\text{Property 66.10}}{=} \frac{1}{n^2} \mathbb{V}\left[\sum_{i=1}^n x_i\right] \\ \frac{1}{n^2} n \mathbb{V}[X] &= \frac{1}{n} \sigma^2 \end{aligned}$$

Proof 68.3: definition 68.11:

$$\begin{aligned} \mathbb{E}[\hat{\sigma}_X^2] &= \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n (x_i^2 - 2x_i \bar{x} + \bar{x}^2)\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - 2n\bar{x} \cdot n\bar{x} + n\bar{x}^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - n\bar{x}^2\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n \mathbb{E}[x_i^2] - n\mathbb{E}[\bar{x}^2] \right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \mu^2) - n\mathbb{E}[\bar{x}^2] \right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \mu^2) - n\left(\frac{1}{n}\sigma^2 + \mu^2\right) \right] \\ &= \frac{1}{n-1} [(n\sigma^2 + n\mu^2) - (\sigma^2 + n\mu^2)] \\ &= \frac{1}{n-1} [n\sigma^2 - \sigma^2] = \frac{1}{n-1} [(n-1)\sigma^2] = \sigma^2 \end{aligned}$$

Stochastic Calculus

Stochastic Processes

Definition 69.1

Random/Stochastic Process

$$\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\}$$

An $(\mathbb{R}^d\text{-valued})$ stochastic process is a collection of $(\mathbb{R}^d\text{-valued})$ random variables X_t on the same probability space $(\Omega, \mathcal{A}, \mathbb{P})$. The index set \mathcal{T} is usually representing time and can be either an interval $[t_1, t_2]$ or a discrete set $\{t_1, t_2, \dots\}$. Therefore, the random process X can be written as a function:

$$X : \mathcal{T} \subseteq \mathbb{R}_+ \times \Omega \mapsto \mathbb{R}^d \iff (t, \omega) \mapsto X(t, \omega) \quad (69.1)$$

Definition 69.2 Sample path/Trajectory/Realization: Is the stochastic/noise signal $r(\cdot, \omega)$ on the index set \mathcal{T} , that we obtain by sampling ω from Ω .

Notation

Even though the r.v. X is a function of two variables, most books omit the argument of the sample space $X(t, \omega) := X(t)$

Corollary 69.1

$$\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\} > 0$$

Strictly Positive Stochastic Processes: A stochastic process $\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\}$ is called strictly positive if it satisfies:

$$X_t > 0 \quad \text{P-a.s.} \quad \forall t \in \mathcal{T} \quad (69.2)$$

Definition 69.3

Random/Stochastic Chain

$$\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\}$$

is a collection of random variables defined on the same probability space $(\Omega, \mathcal{A}, \mathbb{P})$. The random variables are ordered by an associated index set \mathcal{T} and take values in the same mathematical **discrete state space** S , which must be measurable w.r.t. some σ -algebra Σ . Therefore for a given probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and a measurable space (S, Σ) , the random chain X is a collection of S -valued random variables that can be written as:

$$X : \mathcal{T} \times \Omega \mapsto S \iff (t, \omega) \mapsto X(t, \omega) \quad (69.3)$$

Definition 69.4 Index/Parameter Set

$$\mathcal{T}$$

Usually represents time and can be either an interval $[t_1, t_2]$ or a discrete set $\{t_1, t_2, \dots\}$.

Definition 69.5 State Space

$$S$$

Is the range/possible values of the random variables of a stochastic process and must be measurable w.r.t. some σ -algebra Σ .

Sample-vs. State Space

Sample space hints that we are working with probabilities i.e. probability measures will be defined on our sample space.

State space is used in dynamics, it implies that there is a time progression, and that our system will be in different states as time progresses.

Definition 69.6 Sample path/Trajectory/Realization: Is the stochastic/noise signal $r(\cdot, \omega)$ on the index set \mathcal{T} , that we obtain by sampling ω from Ω .

Notation

Even though the r.v. X is a function of two variables, most books omit the argument of the sample space $X(t, \omega) := X(t)$

1. Filtrations

Definition 69.7 Filtration

$$\mathcal{F} = \{\mathcal{F}_t\}_{t \geq 0}$$

A collection $\{\mathcal{F}_t\}_{t \geq 0}$ of sub σ -algebras $\{\mathcal{F}_t\}_{t \geq 0} \in \mathcal{F}$ is called filtration if it is increasing:

$$\mathcal{F}_s \subseteq \mathcal{F}_t \quad \forall s \leq t \quad (69.4)$$

Explanation 69.1 (Definition 69.7). A filtration describes the flow of information i.e. with time we learn more information.

Definition 69.8

Filtered Probability Space

$$(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$$

A probability space $(\Omega, \mathcal{F}, \mathbb{P})$ together with a filtration $\{\mathcal{F}_t\}_{t \geq 0}$ is called a filtered probability space.

Definition 69.9 Adapted Process: A stochastic process $\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\}$ is called adapted to a filtration \mathcal{F} if:

$$X_t \text{ is } \mathcal{F}_t\text{-measurable} \quad \forall t \quad (69.5)$$

That is the value of X_t is observable at time t

Definition 69.10 Predictable Process: A stochastic process $\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\}$ is called predictable w.r.t. a filtration \mathcal{F} if:

$$X_t \text{ is } \mathcal{F}_{t-1}\text{-measurable} \quad \forall t \quad (69.6)$$

That is the value of X_t is known at time $t-1$

Note

The price of a stock will usually be adapted since date k prices are known at date k .

On the other hand the interest rate of a bank account is usually already known at the beginning $k-1$, s.t. the interest rate r_t ought to be \mathcal{F}_{k-1} measurable, i.e. the process $r = (r_k)_{k=1, \dots, T}$ should be predictable.

Corollary 69.2 : The amount of information of an adapted random process is increasing see example 69.1.

2. Martingales

Definition 69.11 Martingales: A stochastic process $X(t)$ is a martingale on a **filtered probability space** $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$ if the following conditions hold:

① Given $s \leq t$ the best prediction of $X(t)$, with a filtration $\{\mathcal{F}_s\}$ is the current expected value:

$$\forall s \leq t \quad \mathbb{E}[X(t) | \mathcal{F}_s] = X(s) \quad \text{a.s.} \quad (69.7)$$

② The expectation is finite:

$$\mathbb{E}[|X(t)|] < \infty \quad \forall t \geq 0 \quad X(t) \text{ is } \{\mathcal{F}_t\}_{t \geq 0} \text{ adapted} \quad (69.8)$$

Interpretation

- For any \mathcal{F}_s -adapted process the best prediction of $X(t)$ is the currently known value $X(s)$ i.e. if $\mathcal{F}_s = \mathcal{F}_{t-1}$ then the best prediction is $X(t-1)$
- A martingale models fair games of limited information.

Definition 69.12 Auto Covariance

$$\gamma(t_2 - t_1)$$

Describes the covariance between two values of a stochastic process $(X_t)_{t \in \mathcal{T}}$ at different time points t_1 and t_2 :

$$\gamma(t_1, t_2) = \text{Cov}[X_{t_1}, X_{t_2}] = \mathbb{E}[(X_{t_1} - \mu_{t_1})(X_{t_2} - \mu_{t_2})] \quad (69.9)$$

For zero time differences $t_1 = t_2$ the autocorrelation functions equals the variance:

$$\gamma(t, t) = \text{Cov}[X_t, X_t] \stackrel{\text{eq. (66.35)}}{=} \text{Var}[X_t] \quad (69.10)$$

Notes

- Hence the autocorrelation describes the correlation of a function or signal with itself at a previous time point.
- Given a random time dependent variable $x(t)$ the autocorrelation function $\gamma(t, t - \tau)$ describes how similar the time translated function $x(t - \tau)$ and the original function $x(t)$ are.
- If there exists some relation between the values of the time series that is non-random, then the autocorrelation is non-zero.
- The auto covariance is maximized/most similar for no translation $\tau = 0$ at all.

Definition 69.13 Auto Correlation

$$\rho(t_2 - t_1)$$

Is the scaled version of the auto-covariance $\gamma(t_2 - t_1)$:

$$\rho(t_2 - t_1) = \text{Corr}[X_{t_1}, X_{t_2}] = \frac{\text{Cov}[X_{t_1}, X_{t_2}]}{\sigma_{X_{t_1}} \sigma_{X_{t_2}}} = \frac{\mathbb{E}[(X_{t_1} - \mu_{t_1})(X_{t_2} - \mu_{t_2})]}{\sigma_{X_{t_1}} \sigma_{X_{t_2}}} \quad (69.11)$$

3. Different kinds of Processes

1. Markov Process

Definition 69.14 Markov Process: A continuous-time stochastic process $X(t), t \in T$, is called a Markov process if for any finite parameter set $\{t_i : t_i < t_{i+1}\} \in T$ it holds:

$$\mathbb{P}(X(t_{n+1}) \in B | X(t_1), \dots, X(t_n)) = \mathbb{P}(X(t_{n+1}) \in B | X(t_n))$$

it thus follows for the **transition probability** – the probability of $X(t)$ lying in the set B at time t , given the value x of the process at time s :

$$\mathbb{P}(s, x, t, B) = \mathbb{P}(X(t) \in B | X(s) = x) \quad 0 \leq s < t \quad (69.12)$$

Interpretation

In order to predict the future only the current/last value counts.

Corollary 69.3 Transition Density: The transition probability of a continuous distribution p can be calculated via:

$$\mathbb{P}(s, x, t, B) = \int_B p(s, x, t, y) dy \quad (69.13)$$

2. Gaussian Process

Definition 69.15 Gaussian Process: Is a stochastic process $X(t)$ where the random variables follow a Gaussian distribution:

$$X(t) \sim \mathcal{N}(\mu(t), \sigma^2(t)) \quad \forall t \in T \quad (69.14)$$

3. Diffusions

Definition 69.16 Diffusion:

[proof 69.1],[proof 69.2]

Is a Markov Process for which it holds that:

$$\mu(t, X(t)) = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \mathbb{E}[X(t + \Delta t) - X(t) | X(t)] \quad (69.15)$$

$$\sigma^2(t, X(t)) = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \mathbb{E}[(X(t + \Delta t) - X(t))^2 | X(t)] \quad (69.16)$$

- $\mu(t, X(t))$ is called **drift**
- $\sigma^2(t, X(t))$ is called **diffusion coefficient**

Interpretation

There exist not discontinuities for the trajectories.

4. Brownian Motion/Wiener Process

Definition 69.17

d-dim standard Brownian Motion/Wiener Process:

Is an \mathbb{R}^d valued stochastic process $(W_t)_{t \in T}$ starting at $x_0 \in \mathbb{R}^d$ that satisfies:

① **Normal Independent Increments:** the increments are normally distributed independent random variables:

$$W(t_i) - W(t_{i-1}) \sim \mathcal{N}(0, (t_i - t_{i-1}) \mathbb{I}_{d \times d}) \quad \forall i \in \{1, \dots, T\} \quad (69.17)$$

② **Stationary increments:**

$W(t + \Delta t) - W(t)$ is independent of $t \in T$

③ **Continuity:** for a.e. $\omega \in \Omega$, the function $t \mapsto W_t(\omega)$ is continuous

$$\lim_{t \rightarrow 0} \frac{\mathbb{P}(|W(t + \Delta t) - W(t)| \geq \delta)}{\Delta t} = 0 \quad \forall \delta > 0 \quad (69.18)$$

④ **Start**

$$W(0) := W_0 = 0 \quad \text{a.s.} \quad (69.19)$$

Notation

- In many source the Brownian motion is a synonym for the standard Brownian Motion and it is the same as the Wiener process.

- However in some sources the Wiener process is the standard Brownian Motion, while the Brownian motion denotes a general form $\alpha W(t) + \beta$.

Corollary 69.4 $W_t \sim \mathcal{N}(0, \sigma)$ [proof 69.4],[proof 69.5]: The random variable W_t follows the $\mathcal{N}(0, \sigma)$ law

$$\mathbb{E}[W(t)] = \mu = 0 \quad (69.20)$$

$$\mathbb{V}[W(t)] = \mathbb{E}[W^2(t)] = \sigma^2 = t \quad (69.21)$$

3.4.1. Properties of the Wiener Process

Property 69.1 Non-Differentiable Trajectories:

The sample paths of a Brownian motion are not differentiable:

$$\begin{aligned} \frac{dW(t)}{dt} &= \lim_{\Delta t \rightarrow 0} \mathbb{E}\left[\left(\frac{(W(t + \Delta t) - W(t))}{\Delta t}\right)^2\right] \\ &= \lim_{\Delta t \rightarrow 0} \frac{\mathbb{E}[W(t + \Delta t) - W(t)]}{\Delta t} = \lim_{\Delta t \rightarrow 0} \frac{\sigma^2}{\Delta t} = \infty \end{aligned}$$

result cannot use normal calculus anymore
solution Ito Calculus see section 70.

Property 69.2 Auto covariance Function:

The auto-covariance $\mathbb{E}[(W(t) - \mu)(W(t') - \mu)]$ for a Wiener process

$$\mathbb{E}[(W(t) - \mu)(W(t') - \mu)] = \min(t, t') \quad (69.22)$$

Property 69.3: A standard Brownian motion is a Quadratic Variation

Definition 69.18 Total Variation: The total variation of a function $f : [a, b] \subset \mathbb{R} \mapsto \mathbb{R}$ is defined as:

$$LV[a, b](f) = \sup_{\Pi \in \mathcal{S}} \sum_{i=0}^{n_{\Pi}-1} |f(x_{i+1}) - f(x_i)| \quad (69.23)$$

$$\mathcal{S} = \{\Pi\{x_0, \dots, x_{n_{\Pi}}\} : \Pi \text{ is a partition of } [a, b]\}$$

it is a measure of the (one dimensional) length of a function w.r.t. to the y-axis, when moving along the function.

Hence it is a measure of the variation of a function w.r.t. to the y-axis.

Definition 69.19

Total Quadratic Variation/“sum of squares”:

The total quadratic variation of a function $f : [a, b] \subset \mathbb{R} \mapsto \mathbb{R}$ is defined as:

$$QV[a, b](f) = \sup_{\Pi \in \mathcal{S}} \sum_{i=0}^{n_{\Pi}-1} |f(x_{i+1}) - f(x_i)|^2 \quad (69.24)$$

$$\mathcal{S} = \{\Pi\{x_0, \dots, x_{n_{\Pi}}\} : \Pi \text{ is a partition of } [a, b]\}$$

Corollary 69.5 Bounded (quadratic) Variation:

The (quadratic) variation of a function is bounded if it is finite:

$$\exists M \in \mathbb{R}_+ : LV[a, b](f) \leq M \quad (QV[a, b](f) \leq M) \quad \forall \Pi \in \mathcal{S} \quad (69.25)$$

Theorem 69.1 Variation of Wiener Process: Almost surely the total variation of a Brownian motion over a interval $[0, T]$ is infinite:

$$\mathbb{P}(\omega : LV(W(\omega)) < \infty) = 0 \quad (69.26)$$

Theorem 69.2

Quadratic Variation of standard Brownian Motion:

The quadratic variation of a standard Brownian motion over $[0, T]$ is finite:

$$\lim_{N \rightarrow \infty} \sum_{k=1}^N \left[W\left(\frac{k}{N}\right) - W\left(\frac{(k-1)}{N}\right) \right]^2 = T \quad (69.27)$$

with probability 1

Corollary 69.6 : theorem 69.2 can also be written as:

$$(dW(t))^2 = dt \quad (69.28)$$

3.4.2. Lévy's Characterization of BM

Theorem 69.3 [proof 69.7],[proof 69.8] **d-dim standard BM/Wiener Process by Paul Lévy:**

An R^d valued adapted stochastic process^[def. 69.1, 69.7] $(W_t)_{t \in T}$ with the filtration $\{\mathcal{F}_t\}_{t \in \mathbb{R}_+}$, that satisfies:

① Start

$$W(0) := W_0 = 0 \quad a.s. \quad (69.29)$$

② **Continuous Martingale:** W_t is an a.s. continuous martingale^[def. 69.11] w.r.t. the filtration $(\mathcal{F}_t)_{t \in T}$ under \mathbb{P} .

③ **Quadratic Variation:**

$$W_t^2 - t \text{ is also a martingale} \iff QV(W_t) = t \quad (69.30)$$

is a standard Brownian motion^[def. 69.24].

Further Stochastic Processes

3.4.3. White Noise

Definition 69.20 Discrete-time white noise: Is a random signal $\{\epsilon_t\}_{t \in T_{\text{discret}}}$ having equal intensity at different frequencies and is defined by:

- Having zero tendencies/expectation (otherwise the signal would not be random):

$$\mathbb{E}[\epsilon * [k]] = 0 \quad \forall k \in T_{\text{discret}} \quad (69.31)$$

- Zero autocorrelation^[def. 69.13] γ i.e. the signals of different times are in-way correlated:

$$\begin{aligned} \gamma(\epsilon * [k], \epsilon * [k+n]) &= \mathbb{E}[\epsilon * [k]\epsilon * [k+n]^T] \\ &= \mathbb{V}[\epsilon * [k]]\delta_{\text{discret}[n]} \end{aligned} \quad \forall k, n \in T_{\text{discret}} \quad (69.32)$$

With $\delta_{\text{discret}[n]} := \begin{cases} 1 & \text{if } n = 0 \\ 0 & \text{else} \end{cases}$

See proofs

Definition 69.21 Continuous-time white noise: Is a random signal $\{\epsilon_t\}_{t \in T_{\text{continuous}}}$ having equal intensity at different frequencies and is defined by:

- Having zero tendencies/expectation (otherwise the signal would not be random):

$$\mathbb{E}[\epsilon * (t)] = 0 \quad \forall t \in T_{\text{continuous}} \quad (69.33)$$

- Zero autocorrelation^[def. 69.13] γ i.e. the signals of different times are in-way correlated:

$$\begin{aligned} \gamma(\epsilon * (t), \epsilon * (t + \tau)) &= \mathbb{E}[\epsilon * (t)\epsilon * (t + \tau)^T] \quad (69.34) \\ \text{eq. (66.91)} \quad \mathbb{V}[\epsilon * (t)]\delta(t - \tau) &= \begin{cases} \mathbb{V}[\epsilon * (t)] & \text{if } \tau = 0 \\ 0 & \text{else} \end{cases} \end{aligned} \quad \forall t, \tau \in T_{\text{continuous}} \quad (69.35)$$

Definition 69.22 Homoscedastic Noise: Has constant variability for all observations/time-steps:

$$\begin{aligned} \mathbb{V}[\epsilon_{i,t}] &= \sigma^2 \quad \forall t = 1, \dots, T \\ \forall i = 1, \dots, N \end{aligned} \quad (69.36)$$

Definition 69.23 Heteroscedastic Noise: Is noise whose variability may vary with each observation/time-step:

$$\begin{aligned} \mathbb{V}[\epsilon_{i,t}] &= \sigma(i, t)^2 \quad \forall t = 1, \dots, T \\ \forall i = 1, \dots, N \end{aligned} \quad (69.37)$$

3.4.4. Generalized Brownian Motion

Definition 69.24 Brownian Motion: Let $\{W_t\}_{t \in \mathbb{R}_+}$ be a standard Brownian motion^[def. 69.17], and define:

$$X_t = \mu t + \sigma W_t \quad t \in \mathbb{R}_+ \quad \mu \in \mathbb{R} : \text{drift parameter} \quad \sigma \in \mathbb{R}_+ : \text{scale parameter} \quad (69.38)$$

then $\{X_t\}_{t \in \mathbb{R}_+}$ is normally distributed with mean μt and variance $t\sigma^2$. $X_t \sim \mathcal{N}(\mu t, \sigma^2 t)$.

Theorem 69.4 Normally Distributed Increments:

If $W(T)$ is a Brownian motion, then $W(t) - W(0)$ is a normal random variable with mean μt and variance $\sigma^2 t$, where $\mu, \sigma \in \mathbb{R}$. From this it follows that $W(t)$ is distributed as:

$$f_{W(t)}(x) \sim \mathcal{N}(\mu t, \sigma^2 t) = \frac{1}{\sqrt{2\pi\sigma^2 t}} \exp\left\{-\frac{(x - \mu t)^2}{2\sigma^2 t}\right\} \quad (69.39)$$

Corollary 69.7: More generally we may define the process: $t \mapsto f(t) + \sigma W_t$ which corresponds to a noisy version of f .

Corollary 69.8

Brownian Motion as a Solution of an SDE: A stochastic process X_t follows a BM with drift μ and scale σ if it satisfies the following SDE:

$$dX(t) = \mu dt + \sigma dW(t) \quad (69.41)$$

$$X(0) = 0 \quad (69.42)$$

3.4.5. Geometric Brownian Motion (GBM)

For many processes $X(t)$ it holds that:

- there exists an (exponential) growth
- that the values may not be negative $X(t) \in \mathbb{R}_+$

Definition 69.25 Geometric Brownian Motion:

Let $\{W_t\}_{t \in \mathbb{R}_+}$ be a standard Brownian motion^[def. 69.17] the stochastic process $S_t^1 \triangleq S^1(t)$ with drift parameter μ and scale σ satisfying the SDE:

$$\begin{aligned} dS_t^1 &= S_t^1 (\mu dt + \sigma dW_t) \\ &= \mu S_t^1 dt + \sigma S_t^1 dW_t \end{aligned} \quad (69.43)$$

is called geometric Brownian motion and is given by:

$$S_t^1 = S_0 \exp\left(\sigma W_t + \left(\mu - \frac{1}{2}\sigma^2\right)t\right) \quad t \in \mathbb{R}_+ \quad (69.44)$$

Corollary 69.9 Log-normal Returns:

For a geometric BM we obtain log-normal returns:

$$\ln\left(\frac{S_t}{S_0}\right) = \bar{\mu}t + \sigma W(t) \iff \bar{\mu}t + \sigma W(t) \sim \mathcal{N}(\mu t, \sigma^2 t)$$

with $\bar{\mu} := \mu - \frac{1}{2}\sigma^2$ (69.45)

3.4.6. Locally Brownian Motion

Definition 69.26 Locally Brownian Motion:

Let $\{W_t\}_{t \in \mathbb{R}_+}$ be a standard Brownian motion^[def. 69.17] a local Brownian motion is a stochastic process $X(t)$ that satisfies the SDE:

$$dX(t) = \mu(X(t), t) dt + \sigma(X(t), t) dW(t) \quad (69.46)$$

Note

A local Brownian motion is an generalization of a geometric Brownian motion.

3.4.7. Ornstein-Uhlenbeck Process

Definition 69.27 Ornstein-Uhlenbeck Process:

Let $\{W_t\}_{t \in \mathbb{R}_+}$ be a standard Brownian motion^[def. 69.17] a Ornstein-Uhlenbeck Process or exponentially correlated noise is a stochastic process $X(t)$ that satisfies the SDE:

$$dX(t) = -aX(t)dt + b\sigma dW(t) \quad a > 0 \quad (69.47)$$

5. Poisson Processes

Definition 69.28 Rare/Extreme Events:

Are events that lead to discontinuous in stochastic processes.

Problem

A Brownian motion is not sufficient as model in order to describe extreme events s.a. crashes in financial market time series. Need a model that can describe such discontinuities/jumps.

Definition 69.29 Poisson Process:

A Poisson Process with rate $\lambda \in \mathbb{R}_{\geq 0}$ is a collection of random variables $X(t)$, $t \in [0, \infty)$ defined on a probability space $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$, having a discrete state space $N = \{0, 1, 2, \dots\}$ and satisfies:

- $X_0 = 0$
- The increments follow a Poisson distribution^[def. 66.27]:

$$\mathbb{P}((X_t - X_s) = k) = \frac{\lambda(t-s)}{k!} e^{-\lambda(t-s)} \quad 0 \leq s < t < \infty \quad \forall k \in \mathbb{N}$$

- No correlation of (non-overlapping) increments:

$$\forall t_0 < t_1 < \dots < t_n : \text{the increments are independent} \quad X_{t_1} - X_{t_0}, X_{t_2} - X_{t_1}, \dots, X_{t_n} - X_{t_{n-1}} \quad (69.48)$$

Interpretation

A Poisson Process is a continuous-time process with discrete, positive realizations in $\mathbb{N}_{\geq 0}$

Corollary 69.10 Probability of events: Using Taylor in order to expand the Poisson distribution one obtains:

$$\mathbb{P}(X_{(t+\Delta t)} - X_t \neq 0) = \lambda \Delta t + o(\Delta t^2) \quad t \text{ small i.e. } t \rightarrow 0 \quad (69.49)$$

- Thus the probability of an event happening during Δt is proportional to time period and the rate λ
- The probability of two or more events to happen during Δt is of order $o(\Delta t^2)$ and thus extremely small (as Δt is small).

Definition 69.30 Differential of a Poisson Process:

The differential of a Poisson Process is defined as;

$$dX_t = \lim_{\Delta t \rightarrow dt} (X_{(t+\Delta t)} - X_t) \quad (69.50)$$

Property 69.4 Probability of Events for differential:

With the definition of the differential and using the previous results from the Taylor expansion it follows:

$$\mathbb{P}(dX_t = 0) = 1 - \lambda \quad (69.51)$$

$$\mathbb{P}(|dX_t| = 1) = \lambda \quad (69.52)$$

Proofs

Proof 69.1: eq. (69.15):

Let by δ denote the displacement of a particle at each step, and assume that the particles start at the center i.e. $x(0) = 0$, then we have:

$$\begin{aligned} \mathbb{E}[x(n)] &= \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N x_i(n)\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1) \pm \delta] \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1)] \\ &\stackrel{\text{ind.}}{=} \mathbb{E}[x_{n-1}] = \dots = \mathbb{E}[x(0)] = 0 \end{aligned}$$

Thus in expectation the particles goes nowhere.

Proof 69.2: eq. (69.16):

Let by δ denote the displacement of a particle at each step, and assume that the particles start at the center i.e. $x(0) = 0$, then we have:

$$\begin{aligned} \mathbb{E}[x(n)^2] &= \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N x_i(n)^2\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1) \pm \delta]^2 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1)^2 \pm 2\delta x_i(n-1) + \delta^2] \\ &\stackrel{\text{ind.}}{=} \mathbb{E}[x_{n-1}^2] + \delta^2 = \mathbb{E}[x_{n-2}^2] + 2\delta^2 = \dots \\ &= \mathbb{E}[x(0)] + n\delta^2 = n\delta^2 \end{aligned}$$

as $n = \frac{\text{time}}{\text{step-size}} = \frac{t}{\Delta x}$ it follows:

$$\sigma^2 = \mathbb{E}[x^2(n)] - \mathbb{E}[x(n)]^2 = \mathbb{E}[x^2(n)] = \frac{\delta^2}{\Delta x} t \quad (69.53)$$

Thus in expectation the particles goes nowhere.

Proof 69.3: eq. (69.34):

$$\begin{aligned} \gamma(\epsilon * [k], \epsilon * [k+n]) &= \mathbb{Cov}[\epsilon * [k], \epsilon * [k+1]] \\ &= \mathbb{E}[(\epsilon * [k] - \mathbb{E}[\epsilon * [k]])(\epsilon * [k+n] - \mathbb{E}[\epsilon * [k+n]])^T] \end{aligned}$$

$$\stackrel{\text{eq. (69.31)}}{=} \mathbb{E}[(\epsilon * [k])(\epsilon * [k+n])]$$

Proof 69.4: [cor. 69.4]:

Since $B_t - B_s$ is the increment over the interval $[s, t]$, it is the same in distribution as the increment over the interval $[s-s, t-s] = [0, t-s]$

Thus $B_t - B_s \sim B_{t-s} - B_0$

but as B_0 is a.s. zero by definition eq. (69.19) it follows: $B_t - B_s \sim B_{t-s}$ $B_{t-s} \sim \mathcal{N}(0, t-s)$

Proof 69.5: [cor. 69.4]:

$$W(t) = W(t) - \underbrace{W(0)}_{=0} \sim \mathcal{N}(0, t)$$

$$\Rightarrow \mathbb{E}[X] = 0 \quad \mathbb{V}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = t$$

Proof 69.6: theorem 69.2:

$$\begin{aligned} \sum_{k=0}^{N-1} [W(t_k) - W(t_{k-1})]^2 &\quad t_k = \frac{k}{N} T \\ &= \sum_{k=0}^{N-1} X_k^2 \quad X_k \sim \mathcal{N}\left(0, \frac{T}{N}\right) \\ &= \sum_{k=0}^{N-1} Y_k = n \left(\frac{1}{n} \sum_{k=0}^{N-1} Y_k \right) \quad \mathbb{E}[Y_k] = \frac{T}{N} \\ &\stackrel{\text{S.L.L.N}}{=} \frac{T}{n} = T \end{aligned}$$

Proof 69.7: theorem 69.3 ②:

- first we need to show eq. (69.7): $\mathbb{E}[W_t | \mathcal{F}_s] = W_s$

Due to the fact that W_t is \mathcal{F}_t measurable i.e. $W_t \in \mathcal{F}_t$ we know that:

$$\mathbb{E}[W_t | \mathcal{F}_t] = W_t \quad (69.54)$$

$$\begin{aligned} \mathbb{E}[W_t | \mathcal{F}_s] &= \mathbb{E}[W_t - W_s + W_s | \mathcal{F}_s] \\ &= \mathbb{E}[W_t - W_s | \mathcal{F}_s] + \mathbb{E}[W_s | \mathcal{F}_s] \\ &\stackrel{\text{eq. (69.54)}}{=} \mathbb{E}[W_t - W_s] + W_s \\ &\stackrel{\text{W}_t - W_s \sim \mathcal{N}(0, t-s)}{=} W_s \end{aligned}$$

- second we need to show eq. (69.8): $\mathbb{E}[|X(t)|] < \infty$

$$\mathbb{E}[|W(t)|]^2 \stackrel{?}{\leq} \mathbb{E}[|W(t)|^2] = \mathbb{E}[W^2(t)] = t < \infty$$

Proof 69.8: theorem 69.3 ③: $W_t^2 - t$ is a martingale?

Using the binomial formula we can write and adding $W_s - W_s$:

$$W_t^2 = (W_t - W_s)^2 + 2W_s(W_t - W_s) + W_s^2$$

using the expectation:

$$\begin{aligned} \mathbb{E}[W_t^2 | \mathcal{F}_s] &= \mathbb{E}[(W_t - W_s)^2 | \mathcal{F}_s] + \mathbb{E}[2W_s(W_t - W_s) | \mathcal{F}_s] \\ &\quad + \mathbb{E}[W_s^2 | \mathcal{F}_s] \end{aligned}$$

$$\stackrel{\text{eq. (69.54)}}{=} \mathbb{E}[(W_t - W_s)^2] + 2W_s \mathbb{E}[(W_t - W_s)] + W_s^2$$

$$\stackrel{\text{eq. (69.21)}}{=} \mathbb{V}[W_t - W_s] + 0 + W_s^2$$

$$t - s + W_s^2$$

from this it follows that:

$$\mathbb{E}[W_t^2 - t | \mathcal{F}_s] = W_s^2 - s \quad (69.55)$$

Examples

Example 69.1 :

Suppose we have a sample space of four elements:
 $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$. At time zero, we do not have any information about which ω has been chosen. At time $T/2$ we know whether we have $\{\omega_1, \omega_2\}$ or $\{\omega_3, \omega_4\}$. At time T , we have full information.

$$\mathcal{F} = \begin{cases} \{\emptyset, \Omega\} & t \in [0, T/2) \\ \{\emptyset, \{\omega_1, \omega_2\}, \{\omega_3, \omega_4\}, \Omega\} & t \in [T/2, T) \\ \mathcal{F}_{\max} = 2^{\Omega} & t = T \end{cases} \quad (69.56)$$

Thus, \mathcal{F}_0 represents initial information whereas \mathcal{F}_{∞} represents full information (all we will ever know). Hence, a stochastic process is said to be defined on a filtered probability space $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$.

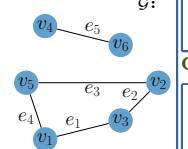
Ito Calculus

[\[Return to Top\]](#)

Graph Theory

Definition 71.1 Graph

A graph \mathcal{G} is a pair $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ of a finite set of vertices \mathcal{V} ^[def. 71.4] and a multi set^[def. 50.3] of edges \mathcal{E} ^[def. 71.10].



Definition 71.2 Order

$n = |\mathcal{V}|$: The order of a graph is the cardinality of its vertex set.

Definition 71.3 Size

$m = |\mathcal{E}|$: The size of a graph is the number of its edges.

Corollary 71.1 n -Graph: Is a graph \mathcal{G} ^[def. 71.1] of order n .

Corollary 71.2 (p, q)-Graph: Is a graph \mathcal{G} ^[def. 71.1] of order p and size q .

Vertices

Definition 71.4 Vertices/Nodes

\mathcal{V} : Is a set of entities of a graph connected and related by edges in some way:

Definition 71.5 Neighborhood $N(v)$: The neighborhood of a vertex $v_i \in \mathcal{V}$ is the set of all adjacent vertices:

$$N(v_i) = \{v_k \in \mathcal{V} : \exists e_k = \{v_i, v_j\} \in \mathcal{E}, \forall v_j \in \mathcal{E}\} \quad (71.1)$$

1.0.1. Adjacency Matrix

Definition 71.6 (unweighted) Adjacency Matrix \mathbf{A} : Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ its adjacency matrix is a square matrix $\mathbf{A} \in \mathbb{N}^{n,n}$ defined as:

$$\mathbf{A}_{i,j} := \begin{cases} 1 & \text{if } \exists e(i, j) \\ 0 & \text{otherwise} \end{cases} \quad (71.2)$$

Definition 71.7 weighted Adjacency Matrix \mathbf{A} : Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ its weighted adjacency matrix is a square matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ defined as:

$$\mathbf{A}_{i,j} := \begin{cases} \theta_{i,j} & \text{if } \exists e(i, j) \\ 0 & \text{otherwise} \end{cases} \quad (71.3)$$

Diagonal Elements

For a graph without self-loops the diagonal elements of the adjacency are all zero.

1.0.2. Degree Matrix

Definition 71.8 Degree of a Vertex

δ : The degree of a vertex v is the cardinality of the neighborhood^[def. 71.5] – the number of adjacent vertices:

$$\deg(v_i) = \delta(v) = |N(v)| = \sum_{j=1}^{j < i} \mathbf{A}_{i,j} \quad (71.4)$$

Definition 71.9 Degree Matrix

\mathbf{D} : Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ its degree matrix is a diagonal matrix $\mathbf{D} \in \mathbb{N}^{n,n}$ defined as:

$$\mathbf{D}_{i,j} := \begin{cases} \deg(v_i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (71.5)$$

Edges

Definition 71.10 Edges

$\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$: Represent some relation between edges^[def. 71.4] and are represented by two-element subset sets of the vertices:

$$e_k = \{v_i, v_j\} \in \mathcal{E} \iff v_i \text{ and } v_j \text{ connected} \quad (71.6)$$

Proposition 71.1 Number of Edges: A graph \mathcal{G} with $n = |\mathcal{V}|$ has between $[0, \frac{1}{2}n(n - 1)]$ edges.

Subgraph

Definition 71.11 Subgraph $\mathcal{H} \subseteq \mathcal{G}$: A graph $\mathcal{H} = (U, F)$ is a subgraph of a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ iff: $U \subseteq \mathcal{V}$ and $F \subseteq \mathcal{E}$ (71.7)

Components

Definition 71.12 Component: A connected component of a graph \mathcal{G} is a connected^[def. 74.1] subgraph^[def. 71.11] of \mathcal{G} that is maximal by inclusion – there exist no larger connected containing subgraphs.

The number of components of a graph \mathcal{G} is defined as $c(\mathcal{G})$.

Walks, Paths and cycles

Definition 71.13 Walk: A walk of a graph \mathcal{G} as a sequence of vertices with corresponding edges:

$$W = \{v_k, v_{k+1}\}_{k=1}^K \in \mathcal{E} \quad (71.8)$$

Definition 71.14 Length of a Walk K : Is the number of edges of that Walk.

Definition 71.15 Path P : Is a walk of a graph \mathcal{G} where all visited vertices are distinct (no-repetitions).

Attention: Some use the terms walk for paths and simple paths for paths.

Definition 71.16 Cycle: Is a path^[def. 71.15] of a graph \mathcal{G} where the last visited vertex is the one from which we started.

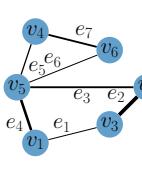
Kinds of Graphs

Weighted Graphs

Definition 72.1 Weighted Graph:

Is a graph \mathcal{G} where edges are associated with a weight:

$$\exists \theta_i := \text{weight}(e_i) \quad \forall e_i \in \mathcal{E}$$

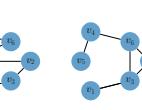


Spanning Graphs

Definition 73.1 Spanning Graph:

Is a subgraph^[def. 71.11] $\mathcal{H} = (U, F)$ of a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ for which it holds:

$$U = \mathcal{V} \quad \text{and} \quad F \subseteq \mathcal{E} \quad (73.1)$$



1. Minimum Spanning Graph

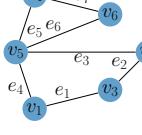
Definition 73.2 Minimum Spanning Graph: Is a spanning graph^[def. 73.1] $\mathcal{H} = (U, F)$ of a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with minimal weights/distance of the edges.

Connected Graphs

Definition 74.1 (Weakly) Connected Graph:

Is a graph \mathcal{G} ^[def. 71.1] where there exists a path between any two vertices:

$$\exists P(v_i, \dots, v_j) \quad \forall v_i, v_j \in \mathcal{V} \quad (74.1)$$



Corollary 74.1 Strongly Connected Graph: A directed graph^[def. 74.3] is called strongly connected if every node is reachable from every other node.

Corollary 74.2 Components of Connected Graphs: A connected graph^[def. 74.1] consists of one component $c(\mathcal{G}) = 1$.

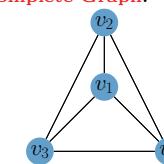
1. Fully Connected/Complete

Definition 74.2 Fully Connected/Complete Graph:

Is a connected graph \mathcal{G} ^[def. 74.1] where each node is connected to every other node.

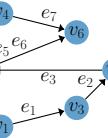
$$\exists e \forall \{v_i, v_j\} \quad \forall v_i, v_j \in \mathcal{V} \quad (74.2)$$

$$|V| = \frac{1}{2} |\mathcal{V}|(|\mathcal{V}| - 1) \quad (74.3)$$



1. Directed Graphs

Definition 74.3 Directed Graph/Digraph (DG):



A directed graph \mathcal{G} is a graph where edges are direct arcs^[def. 74.4].

Definition 74.4 Directed Edges/Arcs: Represent some directional relationship between edges^[def. 71.4] and are represented by ordered two-element subset sets of vertices:

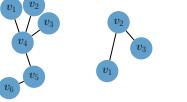
$$e_k = \{v_i, v_j\} \in \mathcal{E} \iff v_i \text{ goes to } v_j \quad (74.4)$$

2. Trees And Forests

2.0.1. Acyclic Graphs

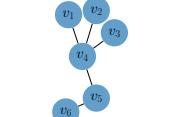
Definition 74.5 Acyclic Graphs: Are graphs^[def. 71.1] where no cycles^[def. 71.16] exist.

Definition 74.6 Forests:



Are acyclic graphs^[def. 74.5]:

Definition 74.7 Trees:



Are acyclic graphs^[def. 74.5] that are connected^[def. 74.1].

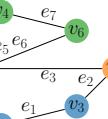
3. Graph Layering

Definition 74.8 Graph Layering:

Given a graph \mathcal{G} a layering of the graph is a partition of its node set

$$\{\mathcal{V}_1, \dots, \mathcal{V}_L\} \subseteq \mathcal{V}$$

$$\text{s.t. } \mathcal{V} = \mathcal{V}_1 \cup \dots \cup \mathcal{V}_L \quad (74.5)$$



4. Bisection Algorithms

1. Local Approaches

2. Global Approaches

4.2.1. Spectral Decomposition

Definition 74.9 Graph Laplacian (Matrix) $L(\mathcal{G})$: Given a graph with n vertices and m edges has a graph laplacian matrix defined as:

$$L = A - D \quad L_{ij} := \begin{cases} -1 & \text{if } i \neq j \text{ and } e_{ij} \in \mathcal{E} \\ 0 & \text{if } i \neq j \text{ and } e_{ij} \notin \mathcal{E} \\ \deg(v_i) & \text{if } i = j \end{cases} \quad (74.6)$$

Corollary 74.3 title: