

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.

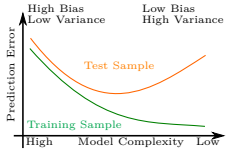
**Definition 1.2 Model Selection:** Is the process of selecting a model  $f$  from a given or chosen class of models  $\mathcal{F}$

**Definition 1.3 Hyperparameter Tuning:** Is the process of choosing the hyperparameters  $\theta$  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  $\Rightarrow$  big test error.



- 1.1. 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  $\mathbb{P}_{\mathcal{X}, \mathcal{Y}}$ .

$$R_P(m) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}}[l(y; m(\mathbf{x}))] = \int_{\mathcal{D}} \mathbb{P}(\mathbf{x}, y) l(y; m(\mathbf{x})) \, d\mathbf{x} \, dy$$
$$= \int_{\mathcal{X}} \int_{\mathcal{Y}} \mathbb{P}(\mathbf{x}, y) l(y; m(\mathbf{x})) \, d\mathbf{x} \, dy$$
$$\stackrel{??}{=} \int_{\mathcal{X}} \int_{\mathcal{Y}} l(y; m(\mathbf{x})) \mathbb{P}(y|\mathbf{x}) \mathbb{P}(\mathbf{x}) \, d\mathbf{x} \, dy \quad (1.1)$$

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  $\mathbf{x}$  but not the distribution of those measurements ( $\mathbf{x} \sim \mathbb{P}_{\mathcal{X}}(\mathbf{x})$ ), we can still calculate the expected risk given/conditioned on the known measurement  $\mathbf{x}$ :

$$\mathcal{R}_P(m, \mathbf{x}) = \int_{\mathcal{Y}} l(y, m(\mathbf{x})) \mathbb{P}(y|\mathbf{x}) \, dy$$

**Corollary 1.1 Note:** [def. 1.6]  $\iff$  [def. 1.7];

$$R_P(m) = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}}[R_P(m, \mathbf{x})] = \int_{\mathcal{X}} \mathbb{P}(\mathbf{x}) R_P(m, \mathbf{x}) \, d\mathbf{x} \quad (1.2)$$

2.1. Expected Risk Minimizer

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

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

3. Empirical Risk

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

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

**Definition 1.9 Empirical Risk:** Is the the average of a loss function of an estimator  $h$  over a finite set of data  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$  drawn from  $\mathbb{P}_{\mathcal{X}, \mathcal{Y}}(\mathbf{x}, y)$ :

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

3.1. Empirical Risk Minimizer

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

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

Questions

- ① How far is the true risk  $\mathcal{R}(m)$  from the empirical risk  $\hat{\mathcal{R}}(m)$ , 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
- $$m^*(\mathbf{x}) \in \arg \min_{m \in \mathcal{F}} \mathcal{R}(m) \quad \text{vs.} \quad \hat{m}(\mathbf{x}) \in \arg \min_{m \in \mathcal{F}} \hat{\mathcal{R}}(m)$$
- We hope that  $\lim_{n \rightarrow \infty} \hat{\mathcal{R}}_n(m) = \mathcal{R}(m)$ .

3.1.1. Squared Loss Expected Squared Risk

**Definition 1.11 Mean Squared Error (MSE):**

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

**Corollary 1.2 title:**

add proof

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

**Definition 1.12 Integrated Means Squared Error (IMSE)/(MISE):** the integrated MSE or *Mean integrated square error* (MISE) is defined as:

$$\text{IMSE} = \int_{\mathcal{X}} \text{MSE}(x) \, d\mathbf{x} = \int_{\mathcal{X}} \mathbb{E}[(\hat{m}(x) - m(x))^2] \, d\mathbf{x} \quad (1.7)$$

Empirical Squared Risk

**Definition 1.13 Mean/Average Squared Prediction Error (MSPE):** the empirical MSE or *Mean/Average Squared Error of Prediction* (MSEP)

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

**Corollary 1.3** [proof 3.2]  
MSEP for new observations: Given a new observation  $x_{\text{new}}$  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  $\sigma$ .

**Definition 1.14** [example 3.9],[proof 3.1]  
**Bayes' optimal predictor for the L2-Loss:**  
**Assuming:** i.i.d. generated data by  $(\mathbf{x}_i, y_i) \sim \mathbb{P}(\mathcal{X}, \mathcal{Y})$ .  
**Considering:** the least squares risk:

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

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

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

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**  $\mathcal{D}/\mathcal{Z}$ :

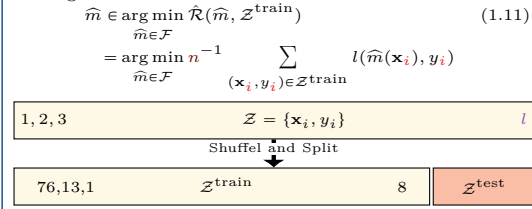
$$\mathcal{Z} = \mathcal{D} := \left\{ \mathbf{z}_j = (\mathbf{x}_j, \mathbf{y}_j) \mid \mathbf{x}_j \in \mathcal{X}, \mathbf{y}_j \in \mathcal{Y} \right\}$$

3.2. 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

$$\mathcal{Z}^{\text{train}} = \left\{ (\mathbf{x}_1^{\text{train}}, y_1^{\text{train}}), \dots, (\mathbf{x}_n^{\text{train}}, y_n^{\text{train}}) \right\}$$

**Definition 1.18 Training Error**  $\hat{\mathcal{R}}(\hat{f}, \mathcal{Z}^{\text{train}})$ : is the model that minimizes the empirical risk [def. 1.10] on the training data [def. 1.17];



3.3. 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.

$$\mathcal{Z}^{\text{test}} = \left\{ (\mathbf{x}_1^{\text{test}}, y_1^{\text{test}}), \dots, (\mathbf{x}_m^{\text{test}}, y_m^{\text{test}}) \right\}$$

**Definition 1.20 Test Error**  $\hat{\mathcal{R}}(\hat{f}, \mathcal{Z}^{\text{test}})$ : 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];

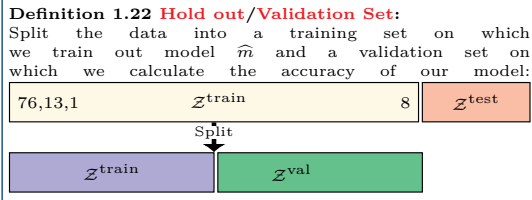
$$\hat{\mathcal{R}}(f, \mathcal{Z}^{\text{test}}) = m^{-1} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{Z}^{\text{test}}} l(\hat{m}(\mathbf{x}_i), y_i) \quad (1.12)$$

3.4. 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.

3.5. Validation Set/Split Once Approach



**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:
$$\hat{m} \in \arg \min_{m \in \mathcal{F}} \hat{\mathcal{R}}(m, \mathcal{Z}^{\text{tr}}) = \arg \min_{m \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n l(y_i, m(\mathbf{x}_i))$$
  - 2: Determine the best parameter  $\theta^*$  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:
$$\hat{\mathcal{R}}(\hat{m}_{\hat{\theta}(\mathcal{Z}^{\text{val}})}(\mathcal{Z}^{\text{tr}}), \mathcal{Z}^{\text{test}})$$

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.

3.6. 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:

$$\hat{m}_{n-1}^{-i} \in \arg \min_{m \in \mathcal{F}} \frac{n-1}{n} \sum_{\substack{j=1 \\ j \neq i}}^n l(y_j, m(x_j)) \quad \forall i \in \{1, \dots, n\}$$
$$\hat{\mathcal{R}}^{\text{LOOCV}} = n^{-1} \sum_{i=1}^n l(y_i, \hat{m}_{n-1}^{-i}(x_i)) \quad (1.13)$$

Pros

- Is basically unbiased estimator, as we use  $n - 1$  training samples.
- Can have a high variance due to highly correlated training sets, as the 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 the 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**:

$$[\hat{m}(x_1) \dots \dots \hat{m}(x_n)]^T = \mathbf{S} \mathbf{Y} \quad (1.14)$$

it holds for the squared loss that:

$$n^{-1} \sum_{i=1}^n (y_i - \hat{m}_{n-1}^{-i}(x_i))^2 = n^{-1} \sum_{i=1}^n \left( \frac{y_i - \hat{m}(x_i)}{1 - \mathbf{S}_{ii}} \right)^2 \quad (1.15)$$

**Definition 1.24 Generalized Cross Validation (GCV):**

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

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

Note

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

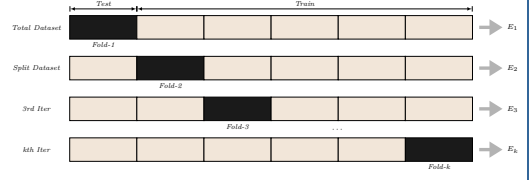
3.7. K-Fold Cross Validation

**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:**

$$\mathcal{Z} = \mathcal{Z}_1 \cup \dots \cup \mathcal{Z}_\nu \cup \dots \cup \mathcal{Z}_K \quad \forall k \in \{1, \dots, K\}$$
$$\widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k} \in \arg \min_{m \in \mathcal{F}} \frac{|\mathcal{Z}_k|}{|\mathcal{Z}|} \sum_{i \in \mathcal{Z} \setminus \mathcal{Z}_k} l(y_i, m(x_i)) \quad (1.17)$$
$$\widehat{\mathcal{R}}^{\text{CV}} = K^{-1} \sum_{k=1}^K |\mathcal{Z}_k|^{-1} \sum_{i \in \mathcal{Z}_k} l\left(y_i, \widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k}(x_i)\right) \quad (1.18)$$

**Note**

A good heuristic for choosing  $K$  is 5, or 10 or:

$$k = \min(\sqrt{n}, 10)$$

**Pros**

- faster then LOOCV.

**Cons**

- runs  $\approx K$  times slower than traing/test-split, as we need to train the model  $K$  times.
- Has higher bias then LOOCV.

There exits systematic tendency to underfit, as each of the  $K$ -fold cross validation models uses only  $n \cdot \frac{K-1}{K}$  training samples

$\Rightarrow$  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.

3.8. Many Random Divisions

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

Generalize LOOCV/ $d$ -fold CV by considering all possible realization?? of  $d$  samples:

$$\mathcal{Z} = \mathcal{Z}_1 \cup \dots \cup \dots \cup \mathcal{Z}_{\binom{n}{d}} \quad \forall k \in \left\{1, \dots, \binom{n}{d}\right\}$$
$$\widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k} \in \arg \min_{m \in \mathcal{F}} \frac{|\mathcal{Z}_k|}{|\mathcal{Z}|} \sum_{i \in \mathcal{Z} \setminus \mathcal{Z}_k} l(y_i, m(x_i)) \quad (1.19)$$
$$\widehat{\mathcal{R}}^{\text{CV}} = \binom{n}{d}^{-1} \sum_{k=1}^{\binom{n}{d}} |\mathcal{Z}_k|^{-1} \sum_{i \in \mathcal{Z}_k} l\left(y_i, \widehat{m}_{n-|\mathcal{Z}_k|}^{-\mathcal{Z}_k}(x_i)\right) \quad (1.20)$$

**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.



# A Statistical Perspective

## 1. Information Theory

### 1.1. 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:

$$I(\mathbf{p}_0 \cdot \mathbf{p}_1 \cdots \mathbf{p}_n) = I(\mathbf{p}_0) + I(\mathbf{p}_1) + \cdots + I(\mathbf{p}_n)$$

$\Rightarrow$  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  $X \sim \mathbf{p}$ :

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

#### 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
$\cong$ 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:

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

### 1.2. 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<sup>[def. 3.2]</sup>:

#### Definition 3.3 Shannon Entropy

[example 3.3]:

Is the expected amount of information of a random variable  $X \sim \mathbf{p}$ :

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

#### Definition 3.4 Differential/Continuous entropy:

Is the continuous version of the Shannon entropy<sup>[def. 3.3]</sup>:

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

#### Notes

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

#### Property 3.1 Non negativity:

Entropy is always non-negative:

$$H(X) \geq 0 \quad \text{if } X \text{ is deterministic} \quad H(X) = 0 \quad (3.4)$$

### 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  $\mathbf{p}_{X|Y}$ . The entropy of  $Y$  conditioned on  $X$  taking a certain value  $x$  is given as:

$$\begin{aligned} H(Y|X=x) &= \mathbb{E}_{Y|X=x} \left[ \log \frac{1}{\mathbf{p}_{Y|X}(Y|X=x)} \right] \\ &= -\mathbb{E}_{Y|X=x} \left[ \log \mathbf{p}_{Y|X}(y|X=x) \right] \end{aligned} \quad (3.5)$$

#### Definition 3.5

proof 3.4

#### Conditional Entropy

$H(Y|X)$ :

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

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

#### Definition 3.6

proof 3.5

#### Chain Rule for Entropy:

$$\begin{aligned} H(Y|X) &= H(X,Y) - H(X) \\ H(X|Y) &= H(X,Y) - H(Y) \end{aligned} \quad (3.7)$$

#### Property 3.2 Monotonicity:

Information/conditioning reduces the entropy  
 $\Rightarrow$  Information never hurts.

$$H(X|Y) \geq H(X) \quad (3.8)$$

#### Corollary 3.2 From eq. (3.17):

$$H(X,Y) \leq H(X) + H(Y) \quad (3.9)$$

### 1.3. Cross Entropy

#### Definition 3.7 Cross Entropy

[proof 3.3]:

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

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

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

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

#### Corollary 3.3 Kullback-Leibler Divergence:

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

### 1.4. Kullback-Leibler (KL) divergence

If we want to measure how different two distributions  $\mathbf{q}$  and  $\mathbf{p}$  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  $\mathbf{p}$  to  $\mathbf{q}$ :** Given two probability distributions  $\mathbf{p}, \mathbf{q}$  of a random variable  $X$ . The Kullback-Leibler divergence is defined to be:

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

and measures how far away a distribution  $\mathbf{q}$  is from a another distribution  $\mathbf{p}$ .

#### Explanation 3.3.

- $\mathbf{p}$  decides where we put the mass if  $\mathbf{p}(x)$  is zero we do not care about  $\mathbf{q}(x)$ .
- $\mathbf{p}(x)/\mathbf{q}(x)$  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.

#### Property 3.3 Non-Symmetric:

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

#### Property 3.4:

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

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

#### Note

The KL-divergence is not a real distance measure as  $\text{KL}(\mathbb{P} \parallel \mathbb{Q}) \neq \text{KL}(\mathbb{Q} \parallel \mathbb{P})$

**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

### 1.5. Jensen-Shanon Divergence

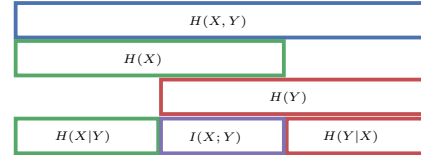
### 1.6. Mutual Information

#### Definition 3.9

example 3.8

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

$$\begin{aligned} I(X;Y) &= H(X) - H(X|Y) = H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X,Y) \\ &= D_{\text{KL}}(\mathbf{p}_{X,Y} \parallel \mathbf{p}_X \mathbf{p}_Y) \end{aligned} \quad (3.17)$$



#### 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??:

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

$$\stackrel{[\text{def. 3.6}]}{\iff} H(z|X) \geq H(z|Y) \quad (3.22)$$

## 2. Proofs

Proof 3.1 Bayes Optimal Predictor<sup>[def. 1.14], :</sup>

$$\begin{aligned} \min_h R(h) &= \min_h \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{p}} [(y - h(\mathbf{x}))^2] \\ &\stackrel{??}{=} \min_h \mathbb{E}_{\mathbf{x} \sim \mathbf{p}_{\mathcal{X}}} [\mathbb{E}_{y \sim \mathbf{p}_{\mathcal{Y}}|\mathcal{X}} [(y - h(\mathbf{x}))^2 | \mathbf{x}]] \\ &\stackrel{?}{=} \mathbb{E}_{\mathbf{x} \sim \mathbf{p}_{\mathcal{X}}} \left[ \underbrace{\min_h \mathbb{E}_{y \sim \mathbf{p}_{\mathcal{Y}}|\mathcal{X}} [(y - h(\mathbf{x}))^2 | \mathbf{x}]}_{\mathcal{R}_{\mathbf{p}}(h, \mathbf{x}) \text{ (def. 1.7)}} \right] \end{aligned}$$

Now lets minimize the conditional executed risk;

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

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

Proof 3.2 Irreducible Error<sup>[cor. 1.3]:</sup>

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

Proof 3.3: Cross Entropy<sup>[def. 3.7]</sup>

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

Notes: ♥

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

Note

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

Proof 3.4: Definition 3.5

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

Proof 3.5: <sup>[def. 3.6]</sup> We start from eq. (3.6):

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

Proof 3.6: example 3.4

$$\begin{aligned} \text{KL}(\mathbf{p}||q) &= \mathbb{E}_{\mathbf{p}} [\log(\mathbf{p}) - \log(q)] \\ &= \mathbb{E}_{\mathbf{p}} \left[ \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2} (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right. \\ &\quad \left. + \frac{1}{2} (\mathbf{x} - \mu_q)^\top \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{p}} \left[ \log \frac{|\Sigma_q|}{|\Sigma_p|} \right] - \frac{1}{2} \mathbb{E}_{\mathbf{p}} \left[ (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right] \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathbf{p}} \left[ (\mathbf{x} - \mu_q)^\top \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] \\ &= \frac{1}{2} \log \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2} \mathbb{E}_{\mathbf{p}} \left[ (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right] \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathbf{p}} \left[ (\mathbf{x} - \mu_q)^\top \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] \\ \underline{\mathbb{E}_{\mathbf{p}}[a]} &\stackrel{\text{tr}(\mathbb{R})}{=} \mathbb{R} \mathbb{E}_{\mathbf{p}} \left[ \text{tr} \left\{ (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right\} \right] \\ &\stackrel{??}{=} \mathbb{E}_{\mathbf{p}} \left[ \text{tr} \left\{ (\mathbf{x} - \mu_p)(\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} \right\} \right] \\ &= \mathbb{E}_{\mathbf{p}} \left[ \text{tr} \left\{ \Sigma_p \Sigma_p^{-1} \right\} \right] \\ &\stackrel{??}{=} \mathbb{E}_{\mathbf{p}} [\text{tr} \{ \mathbf{I}_d \}] = \mathbb{E}_{\mathbf{p}} [d] = d \\ \underline{\mathbb{E}_{\mathbf{p}}[b]} &\stackrel{??}{=} (\mu_p - \mu_q)^\top \Sigma_q^{-1} (\mu_p - \mu_q) + \text{tr} \left\{ \Sigma_q^{-1} \Sigma_p \right\} \end{aligned}$$

### 3. Examples

**Example 3.1 :** Normal distribution has two population parameters: the mean  $\mu$  and the variance  $\sigma^2$ .

**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:**

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

**Example 3.4**

**KL Divergence of Gaussians:**

Given two Gaussian distributions:

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

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

**Example 3.5 KL Divergence of Scalar Gaussians:**

$$\begin{aligned} \theta \sim q(\theta|\lambda) &= \mathcal{N}(\mu_q, \sigma_q^2) \quad \lambda = [\mu_q \quad \sigma_q] \\ \mathbf{p} &= \mathcal{N}(\mu_p, \sigma_p^2) \\ D_{\text{KL}}(\mathbf{p} \parallel q) &= \frac{1}{2} \left( \frac{\sigma_p^2}{\sigma_q^2} (\mu_q - \mu_p)^2 \sigma_q^{-2} - 1 + \log \left( \frac{\sigma_q^2}{\sigma_p^2} \right) \right) \end{aligned}$$

**Example 3.6 KL Divergence of Diag. Gaussians:**

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

**Example 3.7 KL Divergence of Gaussians:**

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

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

**Example 3.8 Gaussian Mutual Information:**

Given  $X \sim \mathcal{N}(\mu, \Sigma)$   $Y = X + \epsilon$   $\epsilon \sim \mathcal{N}(0, \sigma \mathbf{I})$

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

**Example 3.9 Bayes Optimal Predictor and MLE**<sup>[def. 1.14]:</sup>

**Problem:** we do not know the real distribution  $\mathbf{p}_{\mathcal{Y}|\mathcal{X}}(y|\mathbf{x})$ , which we need in order to find the bayes optimal predictor according to eq. (1.10).

**Idea:**

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

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

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

**Idea:** choose parametric form  $\hat{\mathbf{p}}(Y|\mathbf{X}, \theta) = \hat{\mathbf{p}}_\theta(Y|\mathbf{X})$  and then optimize the parameter  $\theta$

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?

$$\mathcal{D} \xrightarrow[\text{Learning method}]{\text{Model Fitting}} (\mathcal{X} \xrightarrow{c} \mathcal{Y}) \xrightarrow[\text{of data } \mathbf{x} \text{ without label}]{\text{Prediction}} \hat{\mathbf{y}}$$

**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.

- **Linear Regression:**  $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$
- **Linear Classification:**  $h(\mathbf{x}) = \text{sing}(\mathbf{w}^\top \mathbf{x})$
- **Kernel Regression:**  $h(\mathbf{x}) = \sum_{i=1}^n \alpha_i \mathbf{k}(\mathbf{x}_i, \mathbf{x})$

- **Neural Networks** (single hidden layer):  
 $h(\mathbf{x}) = \sum_{i=1}^n \mathbf{w}_i' \phi(\mathbf{w}_i^\top \mathbf{x})$

**s.t.** we minimize prediction error/empirical risk <sup>[def. 1.10]</sup>.

**Fundamental assumption**

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

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

**Note**

The distribution  $\mathbf{p}_{\mathcal{X}, \mathcal{Y}}$  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:

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

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  $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathbf{X}$  an estimator  $\hat{\theta} = h(\mathbf{x}_1, \dots, \mathbf{x}_n)$  (3.25)

is a statistic/random variable used to estimate a true (population) parameter  $\theta$  <sup>[def. 3.12]</sup> 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  $\theta$ .

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):**

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

$$\eta = \sum_{j=0}^p \beta_{jm} X_j \tag{3.27}$$

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

Generalized Additive Models (GAMs)

**Definition 3.15** **Generalized Additive Models (GAMs):**

$$sdf \tag{3.29}$$

# 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)<sup>[def. 4.2]</sup>.

## Definition 4.2

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

**Definition 4.3 Coefficients  $\beta$ :** Are the coefficients that we are seeking.

**Definition 4.4 Regression:** Is the process of finding a possible relationship via some coefficients  $\beta$  between *response-variables*  $\mathbf{x}$  and a *predictor-variable(s)*  $\mathbf{y}$  up to some error  $\epsilon$ :  
$$\mathbf{y} = f(\mathbf{x}, \beta) + \epsilon \quad (4.1)$$

## 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  $\beta$  (but not necessarily the data):  
$$\mathbf{y} = \beta^T \phi(\mathbf{x}) + \epsilon \quad (4.2)$$

## Linearity

Linearity is w.r.t. the coefficients  $\beta_j$ .  
Thus a model with transformed non-linear predictor<sup>[def. 4.1]</sup> variables is still called *linear*.

## Definition 4.6 Residual $\mathbf{r}$ :

Let us consider  $n$  observations  $\{x_i, y_i\}_{i=1}^n$ . The residual (error) is the deviation of the observed values from the predicted values:  
$$r_i := e_i = \hat{e}_i = y_i - \hat{y}_i = y_i - \hat{\beta}^T \mathbf{x}_i \quad i = 1, \dots, n \quad (4.3)$$

## Simple (linear) regression (SLR)

**Definition 4.7** <sup>[example 4.1]</sup>  
**Simple Linear Regression:** Is a *linear regression*<sup>[def. 4.8]</sup> with only one explanatory variable<sup>[def. 4.1]</sup>:  
$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i = 1, \dots, n \quad (4.4)$$

## Multiple (linear) regression (MLR)

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

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$$
$$= \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \epsilon_i = \beta^T \mathbf{x}_i + \epsilon_i \quad i = 1, \dots, n$$

$$\begin{bmatrix} \mathbf{x} \end{bmatrix} \begin{bmatrix} \beta \end{bmatrix} = \begin{bmatrix} \mathbf{y} \end{bmatrix} \quad \mathbf{y} = \mathbf{X}\beta \quad \begin{array}{l} \mathbf{X} \in \mathbb{R}^{n, (p+1)} \\ \mathbf{y} \in \mathbb{R}^n \\ \beta \in \mathbb{R}^{p+1} \end{array} \quad \text{Design Matrix:} \quad (4.5)$$

add offset inclusion and fix dimension either p+1 or p or something

## 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

**Simple Linear Quadratic Regression:** Is a *linear regression*<sup>[def. 4.8]</sup> with two explanatory variables<sup>[def. 4.1]</sup> 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:

$$\exists \beta : \begin{bmatrix} x_{11}\beta_1 + x_{12}\beta_2 + \dots + x_{1p}\beta_p \\ x_{21}\beta_1 + x_{22}\beta_2 + \dots + x_{2p}\beta_p \\ \vdots \\ x_{n1}\beta_1 + x_{n2}\beta_2 + \dots + x_{np}\beta_p \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad (4.7)$$

## 1. Linear/Ordinary Least Squares (OLS)

**Problem:** for an over determined system  $n > p$  (usually)  $\forall \mathbf{y} \in \mathfrak{Y}(\mathbf{X})$  (in particular given round off errors) s.t. there exists no parameter vector  $\beta$  that solves<sup>[def. 4.8]</sup>.  
**Idea:** try to find the next best solution by minimizing the residual(s)<sup>[def. 4.6]</sup>.

### Definition 4.10 Residual Sum of Squares:

Is the sum of residuals<sup>[def. 4.6]</sup>:  
$$\text{RSS}(\beta) := \sum_{i=1}^n e_i^2 = \sum_{i=1}^n \|y_i - \hat{y}_i\|_2^2 \quad (4.9)$$

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

Minimizes the residual sum of squares:  
$$\hat{\beta} \in \arg \min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 = \arg \min_{\mathbf{u} \in \mathfrak{Y}(\mathbf{X})} \|\mathbf{y} - \mathbf{u}\|_2^2 \quad (4.10)$$
  
$$= \arg \min_{\beta} \|\mathbf{r}\|_2^2 = \sum_{i=1}^n \left( \sum_{j=1}^p x_{ij}\beta_j - y_i \right)^2 = \text{RSS}(\beta)$$

## 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??.

## 2. Maximum Likelihood Estimate

### Ridge MLE

### Proposition 4.1 (Gauss Markov Assumptions) Assumptions for Linear Regression Model:

1. The  $\{\mathbf{x}_i\}_{i=1}^n$  are deterministic and measured without errors.
2. The variance of the error terms is *homoscedastic*??:  
$$\mathbb{V}[\epsilon_i] = \sigma^2 < \infty \quad \forall i \quad (4.11)$$
3. The errors are uncorrelated:  
$$\text{Cov}[\epsilon_i, \epsilon_j] = 0 \quad \forall i \neq j \quad (4.12)$$
4. The errors are jointly normally distributed with mean 0 and constant variance  $\sigma^2$ :  
$$\epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \forall i = 1, \dots, n \iff \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n) \quad (4.13)$$

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

**Thus:**  $\mathbf{Y}|\mathbf{X} \sim \mathcal{N}(\mathbf{X}\beta, \mathbf{I}\sigma^2)$   $\mathbf{Y}_i|\mathbf{X} \sim \mathcal{N}(\mathbf{x}_i^T \beta, \sigma^2)$   
with:  $\theta = (\beta^T \sigma^2)^T \in \mathbb{R}^{p+1}$

$$l_n(\mathbf{y}|\mathbf{X}, \theta) \propto -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^T \mathbf{x}_i)^2 = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2$$
$$\theta^* \in \arg \max_{\theta \in \mathbb{R}^{p+1}} l_n(\mathbf{y}|\mathbf{X}, \theta) = \arg \min_{\theta \in \mathbb{R}^{p+1}} -l_n(\mathbf{y}|\mathbf{X}, \theta) \quad (4.14)$$

## 2.1. The Normal Equation

### Definition 4.13

**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:

$$\begin{bmatrix} \mathbf{x}^T \mathbf{x} \end{bmatrix} \begin{bmatrix} \beta \end{bmatrix} = \begin{bmatrix} \mathbf{x}^T \mathbf{y} \end{bmatrix} \quad \mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \mathbf{Y} \quad \begin{array}{l} \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{p \times p} \\ \beta \in \mathbb{R}^p \\ \mathbf{X}^T \in \mathbb{R}^{p \times n} \\ \mathbf{Y} \in \mathbb{R}^n \end{array} \quad (4.15)$$

## Geometric Interpretation

### Corollary 4.2

**Geometric Interpretation:**  
We want to find  $\arg \min_{\beta \in \mathbb{R}^n} \|\mathbf{X}\beta - \mathbf{y}\|_2^2$  which is equal to finding:  
$$\arg \min_{\hat{\mathbf{y}} \in \{\mathbf{X}\beta : \beta \in \mathbb{R}^n\} = \mathfrak{Y}(\mathbf{X})} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2$$
but this minimum is equal to the orthogonal projection?? of  $\mathbf{y}$  onto  $\mathfrak{Y}(\mathbf{X})$  i.e. the map:  
$$\mathbf{y} \mapsto \hat{\mathbf{y}}$$
is the orthogonal projection of  $\mathbf{y}$  onto  $\mathfrak{Y}(\mathbf{X})$ .

**Corollary 4.3 Orthogonality of residuals** <sup>[proof 4.6]</sup>:  
Corollary 4.2 implies that the residuals are orthogonal w.r.t. to all the column vectors of  $\mathbf{X}$ :  
$$\mathbf{r}^T \mathbf{x}^{(j)} = 0 \quad \forall j = 1, \dots, p \quad (4.16)$$

### 2.1.1. The Least Squares Solution $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

### Proposition 4.2 Least Squares Solution:

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

## Note

$\mathbf{X}^\dagger$  is the Moore-Penrose pseudo-inverse of the matrix  $\mathbf{X}$ .  
add MPPI to linear algebra appendix

### 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.

1. Compute regular matrix (Matrix Product):  
 $\mathbf{C} := \mathbf{X}^T \mathbf{X} \triangleq \mathcal{O}(n \cdot d^2)$ .
2. Compute the r.h.s. vector (Matrix-Vector):  
 $\mathbf{c} := \mathbf{X}^T \mathbf{y} \in \mathbb{R}^d \triangleq \mathcal{O}(nd)$ .
3. Solve s.p.d. LSE via. **Cholesky decomposition**:  
 $\mathbf{C}\mathbf{w} = \mathbf{c} \triangleq \mathcal{O}(d^3)$ .

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

## Note: s.p.d. $\mathbf{C}$ and cholesky decomposition

**Assume:**  $\mathbf{X}$  has a trivial kernel  $\iff \mathbf{X}^T \mathbf{X}$  is invertible.

1. **Symmetric:** a transposed matrix times itself is symmetric  $\Rightarrow \mathbf{C}$  is symmetric.
2. **Positive definite:**  
 $\mathbf{w}^T \mathbf{C} \mathbf{w} = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} = \|\mathbf{X} \mathbf{w}\|^2 > 0 \quad \forall \mathbf{w} \neq 0$   
has trivial kernel  $\nolimits$

## QR Decomposition

### 2.1.3. Simple Linear Regression Solution

### Definition 4.14

**Linear Regression Solution:** <sup>[proof 4.4]</sup>  
$$\hat{\beta} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{\mathbf{X}^\dagger} \mathbf{y} \quad \text{with} \quad \begin{array}{l} \Sigma^2 = (\mathbf{X}^T \mathbf{X})^{-1} \\ \mathbf{P} = \mathbf{X}^T \mathbf{y} \end{array} \quad (4.18)$$
  
 $\Sigma^2$ : Variance-covar. M.  $\mathbf{P}$ : Inp./Oup. Covariance  
**Moore-Penrose pseudo-inverse:**  $\mathbf{X}^\dagger$  with  $\mathbf{X}^\dagger \mathbf{X} = \mathbf{I}$  (4.19)

### 2.1.4. Making Predictions

**Definition 4.15**  $\mathbf{P}/\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T : \mathbf{y} \mapsto \hat{\mathbf{y}}$   
**Hat/Projection Matrix:** Is the matrix that projects the  $\mathbf{y}$  onto the  $\hat{\mathbf{y}}$ :  
$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} =: \mathbf{P} \mathbf{y} \quad (4.20)$$

**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:

$$\text{tr}(\mathbf{P}) = \text{tr}(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \text{tr}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}) = \text{tr}(\mathbf{I}_{p \times p}) = p$$

**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  $\mathbf{A} \in \mathbb{R}^{p \times p}$ ,  $p \geq p$  then it holds that:  
$$\mathbf{N}(\mathbf{A}) = \mathbf{N}(\mathbf{A}^T \mathbf{A}) \quad \mathfrak{Y}(\mathbf{A}^T) = \mathfrak{Y}(\mathbf{A}^T \mathbf{A}) \quad (4.22)$$

add rest from formulary

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

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (4.23)$$
$$\iff \mathbf{N}(\mathbf{X}) = \{0\} \iff \text{rank}(\mathbf{X}) = p \quad p \geq p \quad (4.24)$$

## 2.2. Moments and Distributions

**Property 4.4 Moments of  $\hat{\beta}$**  <sup>[proof 4.7]</sup>:  
$$\mathbb{E}[\hat{\beta}] = \beta \quad \mathbb{V}[\hat{\beta}] = \text{Cov}[\hat{\beta}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \quad (4.25)$$

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

**Property 4.5 Moments of  $\hat{\mathbf{y}}$**  <sup>[proof 4.9]</sup>:  
$$\mathbb{E}[\hat{\mathbf{y}}] = \mathbb{E}[\mathbf{y}] = \mathbf{X}\beta \quad \mathbb{V}[\hat{\mathbf{y}}] = \text{Cov}[\hat{\mathbf{y}}] = \sigma^2 \mathbf{P} \quad (4.27)$$

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

### Property 4.6 Moments of $\mathbf{r}$ :

$$\mathbb{E}[\mathbf{r}] = 0 \quad \text{Cov}[\mathbf{r}] = \sigma^2 (\mathbf{I} - \mathbf{P}) \quad (4.29)$$
$$\mathbf{r} \sim \mathcal{N}_n(0, \sigma^2 (\mathbf{I} - \mathbf{P})) \quad (4.30)$$

### Property 4.7 Moments of $\hat{\sigma}$ :

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

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

## Note

The standard deviation  $\sigma^2$  is given by  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . However we may not know  $\sigma^2$ , 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  $\sigma$ :

2.2.1. The Gaus Markov Theorem

**Theorem 4.3 Gauss–Markov theorem** [proof 4.10]:  
The BLUE of the  $\beta$  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).

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

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

(4.33)



add delete next section

### 3. MLE with linear Model & Gaussian Noise

#### 3.1. MLE for conditional linear Gaussians

**Questions:** what is  $\mathbb{P}(Y|X)$  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.

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

$\mathbf{X}$ : set of explicative variables.  $\epsilon$ : noise/error term.

**Lemma 4.1 :** The conditional distribution  $D$  of  $Y$  given  $\mathbf{X}$  is equivalent to the unconditional distribution of the noise  $\epsilon$ :

$$\mathbb{P}(Y|\mathbf{X}) \sim D \iff \epsilon \sim D$$

#### 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  $\mathbb{E}[\epsilon] = 0$  and  $y_i = \mathbf{w}^\top \mathbf{x}_i + \epsilon_i$ , as well as ?? it follows:

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

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

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

$$L_n(Y|\mathbf{X}, \theta) = \hat{p}(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \theta)$$

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

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

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

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

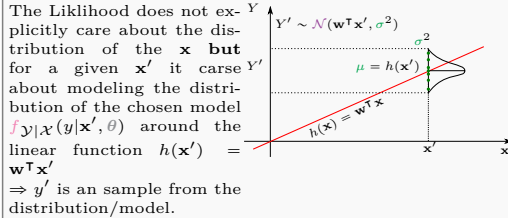
$$\begin{aligned} \frac{\partial \ln(Y|\mathbf{X}, \theta)}{\partial \mathbf{w}} &= \frac{1}{\sigma^2} \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{w}^\top \mathbf{x}_i) = \mathbf{0} \in \mathbb{R}^d \\ &= \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{w} = \sum_{i=1}^n \mathbf{x}_i y_i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

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

#### Note

- The mean  $\mu$  of the normal distribution follows from: 
$$\mathbb{E}[\mathbf{w}^\top \mathbf{x}_i + \epsilon_i] = \mathbb{E}[\mathbf{w}^\top \mathbf{x}_i] + \mathbb{E}[\epsilon_i] = \mathbf{w}^\top \mathbf{x}_i$$

const = 0
- The noise  $\epsilon$  must have zero mean, otherwise it wouldn't be random anymore.
- The optimal function  $h^*(\mathbf{x})$  determines the mean  $\mu$ .
- We can also minimize: 
$$\theta^* = \arg \max_{\theta} \hat{p}(Y|\mathbf{X}, \theta) = \arg \min_{\theta} -\hat{p}(Y|\mathbf{X}, \theta)$$



#### 3.2. Conditional MLE $\hat{=}$ Least Squares

**Assuming** that the noise is i.i.d. Gaussian with constant variance  $\sigma$ , that is  $\theta = (\mathbf{w}^\top \quad \sigma)^T$

and considering the negative log likelihood in order to minimize  $\arg \max_{\alpha} \alpha = -\arg \min_{\alpha}$ :

$$-\ln(\mathbf{w}) = -\prod_{i=1}^n \ln \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \sigma^2) = \frac{n}{2} \ln(2\pi\sigma^2) + \sum_{i=1}^n \frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}$$

$$\arg \max_{\mathbf{w}} \ln(\mathbf{w}) \iff \arg \min_{\mathbf{w}} -\ln(\mathbf{w})$$

$$\arg \min_{\mathbf{w}} \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 = \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \quad (4.35)$$

**Thus** Least squares regression equals Conditional MLE with a linear model + Gaussian noise.

Maximizing Likelihood  $\iff$  Minimizing least squares

**Corollary 4.7 :** The Maximum Likelihood Estimate (MLE) for i.i.d. Gaussian noise (and general models) is given by the squared loss/Least squares solution, assuming that the variance is constant.

#### Heuristics for ??

**Consider** a sample  $\{y_1, \dots, y_n\} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$

$$\frac{\partial \ln(y|\mathbf{x}, \theta)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu) \stackrel{!}{=} 0$$

$$\frac{\partial \ln(y|\mathbf{x}, \theta)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i - \mu)^2 \stackrel{!}{=} 0$$

$$\theta^* = \begin{pmatrix} \mu_*^* \\ \sigma_*^2 \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n y_i \\ \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \end{pmatrix} \quad (4.36)$$

So, the optimal MLE correspond to the empirical mean and the variance.

#### Note

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

#### 3.3. MLE for general conditional Gaussians

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

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

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

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 chose.

#### Ridge Max Prior

##### Prior

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

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

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

eventually add change of variables formula and check if factor is missing

#### Max Prior

$$\begin{aligned} \beta^* &\in \arg \max_{\beta \in \mathbb{R}^{d+1}} \ln(\beta|\Sigma) \\ &= \arg \max_{\beta \in \mathbb{R}^{d+1}} -\frac{1}{2} \ln \det(\Sigma)^{-1} - \frac{d+1}{2} \ln 2\pi - \frac{1}{2} \beta^\top \Sigma^{-1} \beta \\ 0 &\stackrel{!}{=} \frac{\partial}{\partial \beta^*} \ln(\beta^*|\Sigma) = -\frac{\partial}{\partial \beta^*} \beta^* \Sigma^{-1} \beta^* \stackrel{\text{eq. (4.46)}}{=} -2\Sigma^{-1} \beta^* \\ \beta^* &\in \arg \max_{\beta \in \mathbb{R}^{d+1}} \log p(\beta|\Sigma) = \arg \min_{\beta \in \mathbb{R}^{d+1}} -\ln(\beta|\Sigma) = 2\Sigma^{-1} \beta^* \end{aligned} \quad (4.38)$$

#### Log-MAP

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

**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.

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

$\Gamma$ : 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**.

#### Gaussian Prior/Likelihood MAP inference

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

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

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

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  $\lambda$  the bigger can the data faith fullness term be  $\|\mathbf{y} - \mathbf{X}\beta\|^2$ .

#### Note

The intercept  $\beta_0$  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  $\lambda$  increases but will not become zero (as long as  $\lambda \neq \infty$ )!

They are fit in a restricted fashion controlled by the **shrinkage penalty**  $\lambda$ .

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

$\Rightarrow$  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 faith fullness because:

- ①  $\lambda \xrightarrow{\sigma \uparrow} \infty$ :  $\|\beta\|^2$  must be minimized:
  - $\sigma \uparrow$ : model does not need to match data so perfectly as we have more noise in our data/observations  $\iff$  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  $\lambda \xrightarrow{\sigma \downarrow} 0$ :  $\|\mathbf{y} - \mathbf{X}\beta\|^2$  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  $\iff$  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  $\Lambda^{-1} = \mathbf{1} \in \mathbb{R}^{d+1 \times d+1}$
- $\Lambda$  is symmetric and diagonal.
- $(d+1)$  dimension as we included offset into  $\beta$ .

#### Heuristic Map Inference

A really large weight vector  $\beta$  will result in amplifying noise/larger variance/fluctuations  $\triangleq$  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  $\mathbf{X}^T \mathbf{X}$  in eq. (4.40) is not invertible/regular/has not full rank. Then  $(\mathbf{X}^T \mathbf{X} + \Lambda)$  will still be invertible/well posed. This was the original reason for L2/Ridge Regression.

#### MAP $\triangleq$ Ridge

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

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

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

Ridge Regression MAP

Thus if we know our data  $\beta, \sigma$  we can chose  $\lambda$  statistically and do not need cross-validation.

#### Generalization

Regularized estimation can often be understood as MAP inference:

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n l(\mathbf{w}^T \mathbf{x}_i; \mathbf{x}_i, y_i) + C(\mathbf{w}) =$$

$$= \arg \max_{\mathbf{w}} \prod_{i=1}^n \mathbb{P}(\mathbf{w}) \mathbb{P}(y_i|\mathbf{x}_i, \mathbf{w}) = \arg \max_{\mathbf{w}} \mathbb{P}(\mathbf{w}|\text{data})$$

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

#### Priors

### 3.4. Laplace Prior $\triangleq$ Lasso/L1-regularization

#### Intro

- Question:** what if  $d \gg n$  e.g.
- bag of words with  $d = \text{nb. of words} \gg \text{nb. of documents}$ .
  - Genome analysis  $d = \text{nb. of genes} \gg \text{nb. of patients}$ .
- 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. more dim. than observations

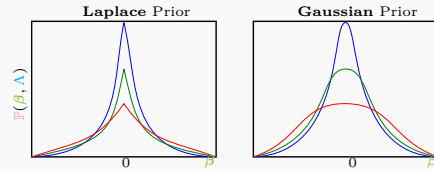
**Idea:** If most of the dimensions are irrelevant for the problem, then we can find a good (**sparse**) solution  $\triangleq$  **feature selection**/dimensionality reduction.

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

$$\mathbb{P}_{\text{Lasso}}(\beta|\Lambda) \stackrel{??}{=} \frac{\Lambda}{2} e^{(-\Lambda|\beta|)} = \prod_{j=1}^d \frac{\lambda_j}{2} e^{-\lambda_j |\beta_j|}$$

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

This leads to a L1 regularized model:



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

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

#### Laplacian MAP Prior Inference

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^d} \left\{ \|\mathbf{y} - (\mathbf{X}\beta + \beta_0)\|^2 + \lambda \|\beta\|_1 \right\}$$

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

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

#### 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 ??) 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.

### 3.5. Sparseness Priors/L0-regularization

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

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

#### 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

Their are two equivalent ways of trading:

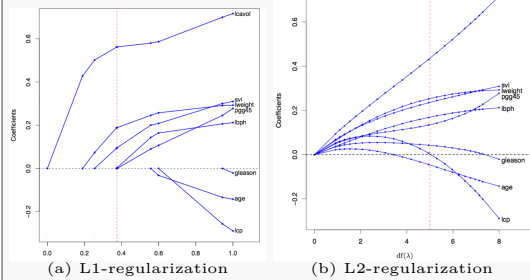
<ul style="list-style-type: none"> <li><math>g(\beta) = \ \mathbf{y} - \mathbf{X}\beta\ ^2</math>: the data term and</li> <li><math>f(\beta)</math>: the Regularizer.</li> </ul> <p><b>Scalarization</b></p> $\min_{\mathbf{w}} f(\beta) + Cg(\beta)$ <p><math>g(\mathbf{w})</math> <math>\begin{cases} C \text{ small} \\ C \text{ large} \end{cases}</math></p>	<p><b>Constraint opt.</b></p> $\min_{\mathbf{w}} g(\beta) \text{ s.t. } f(\mathbf{w}) \leq B$ <p><math>g(\mathbf{w})</math> <math>\begin{cases} B \text{ small} \\ B \text{ large} \end{cases}</math></p>
---	---

#### Note

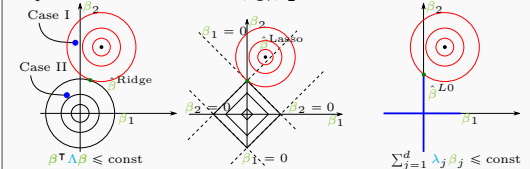
**Scalarization and constrained optimization** gives the same curves  $\iff f, g$  are both convex functions.

**This** is not necessarily for the same values of  $C$  and  $B$  but their exisits 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  $\beta_1, \beta_2$  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  $\triangleq$  corner, due to the fact that corner regions increas 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

#### 3.6. Student's-t likelihood loss function

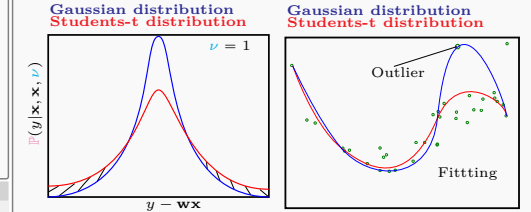
Students-t Distribution:

$$f(\mathbf{y}|\mathbf{x}, \mathbf{w}, \nu, \sigma^2) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\pi\nu\sigma^2}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{(\mathbf{y} - \mathbf{w}^T \mathbf{x})^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}}$$

$\nu$ : 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:**  $\mathbb{P}(|\mathbf{y} - \mathbf{w}^T \mathbf{x}| > t)$  probability of having a outlier/deviation of larger than t, for linear regression.

<b>Students-t</b>	$\mathbb{P}( \mathbf{y} - \mathbf{w}^T \mathbf{x}  > t) = \mathcal{O}(t^{-\alpha})$	$\alpha > 0$
	(Polynomial decay)	
<b>Gaussian</b>	$\mathbb{P}( \mathbf{y} - \mathbf{w}^T \mathbf{x}  > t) = \mathcal{O}(\exp^{-\alpha t})$	$\alpha > 0$
	(Exponential decay)	

$\Rightarrow$  **Students-t** distribution decays less fast then 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:

$$L_n(\mathbf{y}|\mathbf{X}, \theta) \stackrel{\text{i.i.d.}}{=} \prod_{i=1}^n p(\mathbf{x}_i, y_i|\theta) = \prod_{i=1}^n \mathcal{N}(\beta^T \mathbf{x}_i, \sigma^2)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{\sigma^2 2\pi}} \exp\left(-\frac{(y_i - \beta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

$$= (\sigma^2 2\pi)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^T \mathbf{x}_i)^2\right)$$

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

**Proof 4.3 Definition 4.14:**

$$\beta^* \stackrel{\text{i.i.d.}}{=} \arg \min_{\beta \in \mathbb{R}^p} -l_n(\mathbf{y}|\mathbf{X}, \theta)$$

$$= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^T \mathbf{x}_i)^2$$

$$= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$

$$= \arg \min_{\beta \in \mathbb{R}^p} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$

$$\stackrel{\star}{\iff} (-2\mathbf{y}^T \mathbf{X} + 2\mathbf{X}^T \mathbf{X} \beta^*) \stackrel{!}{=} 0$$

$$\Rightarrow \mathbf{X}^T \mathbf{X} \beta^* = \mathbf{X}^T \mathbf{y}$$

Note: ★

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

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

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

Thus

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

combine proofs

Proof 4.4: [def. 4.13]

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

$$0 = \frac{\partial}{\partial \beta} \text{lsq}(\mathbf{X}, \mathbf{y}) = 2\mathbf{X}^\top \mathbf{X}\beta - 2\mathbf{X}^\top \mathbf{y} = 2\mathbf{X}^\top (\mathbf{X}\beta - \mathbf{y})$$

Note

$$\frac{\partial}{\partial \beta} \beta^\top \mathbf{X}^\top (\mathbf{X}\beta) \stackrel{??}{=} (\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \beta = 2\mathbf{X}^\top \mathbf{X}\beta$$

Proof 4.5: Corollary 4.2

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

where  $\mathbf{X} = \{\mathbf{x}_{:,1}, \dots, \mathbf{x}_{:,m}\}$  is the “basis” of the Range space:

$$(\mathbf{X}\beta - \mathbf{y})^\top \mathbf{x}_{:,j} = \mathbf{0} \quad \forall j = 1, \dots, m$$

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

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

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

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

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

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

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  $\mathbf{y}$ :

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

## 5. Examples

Example 4.1 Simple Linear Regression:

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

Example 4.2 Simple Linear Quadratic Regression:

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

# Classification

## 6. Intro

**Definition 4.19 Training Data**  $\mathcal{D}$ :  
 $\mathcal{D} := \{(\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^d, y_i \in \mathcal{Y} := \{c_1, \dots, c_K\}\}$

**Definition 4.20 Classifier**  $c$ :  
Is a mapping that maps the features into classes:  
 $c: \mathcal{X} \rightarrow \mathcal{Y}$  (4.48)

**Definition 4.21 Dichotomy:**  
Given a set  $S = \{s_1, \dots, s_N\}$  a dichotomy is partition of the set  $S$  into two subsets  $A, A^c$  that satisfy:

- Collectively/jointly exhaustiveness:  
 $S = A \cup A^c$  (4.49)
- Mutual exclusivity:  
 $s \in A \implies s \notin A^c \quad \forall s \in S$  (4.50)

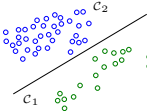
**Explanation 4.1.** Nothing can belong simultaneously to both parts  $A$  and  $A^c$ .

### Types of Classification

**Definition 4.22 Binary Classification:**

Is a classification problem where the labels are binary:

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



### multiclass

### 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 mult-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*?? with one non-zero entry, where the index  $j$  of the non-zero entry indicates the class  $c_j$ :

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

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<sup>[def. 4.26]</sup>:

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

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

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

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

### 7. Binary Classification

$\{-1, 1\}$

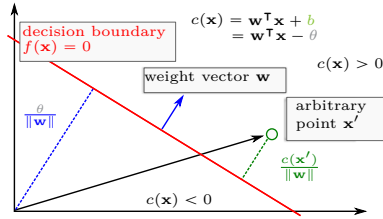
#### 7.1. Linear Classification

**Definition 4.29 Linear Dichotomy:**

**Definition 4.30 Linear Classifier:** A linear classifier is a classifier  $c$  that assigns labels  $\hat{y}$  to samples  $\mathbf{x}_i$  using a linear decision boundary/hyperplane??:

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

**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  $\theta$  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:**

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

**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:

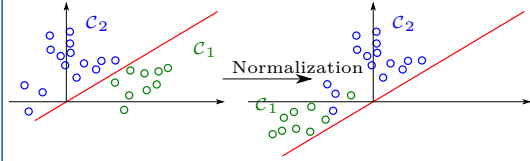
$$\begin{aligned} \mathbf{w} &\leftarrow \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} & \mathbf{x} &\leftarrow \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \\ \implies & \mathbf{w}^\top \mathbf{x} = (\mathbf{w}^\top \ b) \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} = \mathbf{w}^\top \mathbf{x} + b \end{aligned}$$

**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  $y_i$ :

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

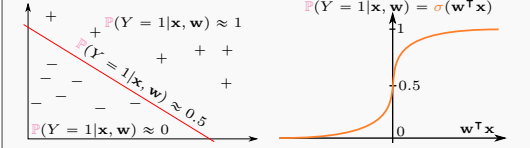
### 8. Logistic Regression

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

**Idea:** in order to classify dichotomies<sup>[def. 4.21]</sup> we use a distribution that maps probabilities to a binary values 0/1  $\Rightarrow$  Bernoulli Distribution??.

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

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



#### 8.1. 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).

$$\sigma(z) = \begin{cases} 0 & -z \text{ large} \\ \frac{1}{2} & \text{if } z \text{ large} \\ 1 & z = 0 \end{cases}$$

#### 8.2. Logistic Regression

**Definition 4.34 Logistic Regression:**

models the likelihood of the output  $y$  as a Bernoulli Distribution??  $y \sim \text{Bern}(\mathbf{p})$ , where the probability  $\mathbf{p}$  is given by the Sigmoid function<sup>[def. 4.33]</sup> of a linear regression:

$$\begin{aligned} \mathbf{p}(y|\mathbf{x}, \mathbf{w}) &= \text{Bern}(\sigma(\mathbf{w}^\top \mathbf{x})) = \begin{cases} \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} & \text{if } y = +1 \\ 1 - \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} & \text{if } y = -1 \end{cases} \\ \stackrel{?? 4.11}{=} & \frac{1}{1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x})} = \sigma(-y \cdot \mathbf{w}^\top \mathbf{x}) \end{aligned} \quad (4.59)$$

#### 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<sup>[def. 6.3]</sup> 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:**

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

### 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<sup>[cor. 4.10]</sup>.

**Definition 4.36 Logistic Loss Gradient** proof 4.13  $\nabla_{\mathbf{w}} l_l(\mathbf{w})$ :

$$\begin{aligned} \nabla_{\mathbf{w}} l_l(\mathbf{w}) &= \mathbb{P}(Y = -y|\mathbf{x}, \mathbf{w}) \cdot (-y\mathbf{x}) \\ &= \frac{1}{1 + \exp(y\mathbf{w}^\top \mathbf{x})} \cdot (-y\mathbf{x}) \end{aligned} \quad (4.62)$$

**Explanation 4.6.**

$$\nabla_{\mathbf{w}} l_l(\mathbf{w}) = \mathbb{P}(Y = -y|\mathbf{x}, \mathbf{w}) \cdot (-y\mathbf{x}) \propto \nabla_{\mathbf{w}} l_H(\mathbf{w})$$

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

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

**Algorithm 4.1 Vanilla SGD for Logistic Regression:**

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

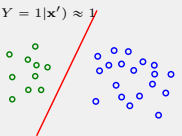
### Making Predictions

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):

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

### Drawback

Logistic regression, does not tell us anything about the likelihood  $\mathbb{P}(Y = 1|\mathbf{x}') \approx 1$  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

### 8.3. Logistic regression and regularization

#### Adding Priors to Logistic Likelihood

• **L2 (Gaussian prior):**

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)) + \lambda \|\mathbf{w}\|_2^2$$

• **L1 (Laplace prior):**

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)) + \lambda \|\mathbf{w}\|_1$$

• **Generalized:**

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)) + \lambda C(\mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \mathbb{P}(\mathbf{w}|\mathbf{X}, Y) \end{aligned}$$

**8.4. SGD for L2-gregularized logistic regression**

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

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

## 9. Proofs

Proof 4.11: <sup>[def. 4.34]</sup> We need to only proof the second expres-  
sion, 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: <sup>[def. 4.35]</sup>

$$l_n(\mathbf{w}) = \arg \max_{\mathbf{w}} \mathbb{P}(y_{1:n}|\mathbf{x}_{1:n}, \mathbf{w}) = \arg \min_{\mathbf{w}} -\log \mathbb{P}(Y|\mathbf{X}, \mathbf{w})$$

$$\stackrel{\text{i.i.d.}}{=} \arg \min_{\mathbf{w}} \sum_{i=1}^n -\log \mathbb{P}(y_i|\mathbf{x}_i, \mathbf{w})$$

$$\stackrel{\text{eq. (4.59)}}{=} -\log \frac{1}{1 + \exp(-y_i \cdot \mathbf{w}^\top \mathbf{x}_i)}$$

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

Proof 4.13: <sup>[def. 4.36]</sup>

$$\nabla_{\mathbf{w}} l_l(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} \log (1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x}))$$

$$\stackrel{\text{C.R.}}{=} \frac{1}{(1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x}))} \frac{\partial}{\partial \mathbf{w}} (1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x}))$$

$$\stackrel{\text{C.R.}}{=} \frac{1}{(1 + \exp(-y \cdot \mathbf{w}^\top \mathbf{x}))} \exp(-y \cdot \mathbf{w}^\top \mathbf{x}) \cdot (-y \mathbf{x})$$

$$= \frac{e^{-z} \cdot (-yx)}{(1 + e^{-z})} = \frac{-yx}{e^z (1 + e^{-z})} = \frac{-yx}{(e^z + e^{-z} + z)}$$

$$= \frac{1}{\exp(y \cdot \mathbf{w}^\top \mathbf{x}) + 1} \cdot (-yx)$$

$$\stackrel{\text{eq. (4.59)}}{=} \hat{p}(Y = -y|\mathbf{x}, \mathbf{w}) \cdot (-y \mathbf{x})$$

# Generalized Linear Models (GLMs)

## 1. Generalized Additive Models (GAMs)

Definition 5.1

Generalized Additive Models (GAM):

Are generalized linear model where the response variable depends linearly on unknown smooth functions  $g_j$  s.t.:

$g_{\text{add}}(\mathbf{x}) = \mu + \sum_{j=1}^p g_j(x_j)$

$g_j : \mathbb{R} \mapsto \mathbb{R}$

$\mathbb{E}[g_j(x_j)] = 0$

$\forall j \in \{1, \dots, p\}$

(5.1)

Pros

- Does not suffer from the curse of dimensionality.

Cons

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

### 1.1. Backfitting

add some point check again Deviance and R squared adj from R output



# Model Parameter Estimation

## 1. Maximum Likelihood Estimation

### 1.1. Likelihood Function

Is a method for estimating the parameters  $\theta$  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  $f_{\mathcal{X}}(\mathbf{x}; \theta)$

**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  $\theta$ :

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

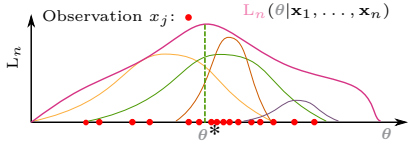


Figure 5: Possible Likelihood function in pink. Overlaid: possible candidate functions for Gaussian model explaining the observations.

#### 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:

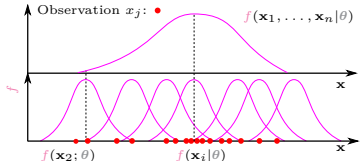


Figure 6: Bottom: probability distributions of the different data points  $\mathbf{x}_i$  given a fixed  $\theta$  for a Gaussian distribution. Top: joint probability distribution of the i.i.d. data points  $\{\mathbf{x}_i\}_{i=1}^n$  given a fixed  $\theta$

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

#### Notation

- The probability density  $f(\mathbf{X}|\theta)$  is considered for a fixed  $\theta$  and thus as a function of the samples.
- The likelihood function on the other hand is considered as a function over parameter values  $\theta$  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  $\theta$  resp.  $\mathbf{X}$  is a parameter and not a random variable.

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

**Definition 6.2 Log Likelihood Function**  $l_n : \Theta \times \mathbb{R}^n \mapsto \mathbb{R}$ :

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

**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*-??:

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

**Definition 6.3 Maximum Likelihood Estimator**  $\theta^*$ : Is the estimator  $\theta^* \in \Theta$  that maximizes the likelihood of the model/predictor:

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

### 1.3. Maximization vs. Minimization

For optimization problems we minimize by convention. The logarithm is a concave function??, 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?? by multiplying it by minus one and consider it as a loss function instead of a likelihood.

**Definition 6.4 Negative Log-likelihood**  $-l_n(\theta|\mathbf{X})$ :

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

### 1.4. 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 an 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 ^ 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 likelihood of a model with parameters  $\theta$  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.

#### 2.1. Prior Distribution

**Definition 6.6 Prior (Distribution)**  $\pi(\theta) = p(\theta)$ :

**Assumes:** that the model parameters  $\theta$  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  $\theta$  and model as synonymous, as the model is fully described by its population parameters<sup>[def. 3.12]</sup>  $\theta$ .

**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  $\pi(\theta|\lambda)$  depends on a set of parameters  $\lambda$ . 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

### 2.2. 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*  $\theta$  a posterior/after observing evidence of a random quantity  $\mathbf{Z}$  that is in a relation with  $\theta$ :

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

**Definition 6.9** [proof 22.1]

**Posterior Distribution and Bayes Theorem:** Using Bayes ?? we can write the posterior distribution as a product of the *likelihood*<sup>[def. 6.1]</sup> weighted with our *prior*<sup>[def. 6.6]</sup> and normalized by the *evidence*  $\mathbf{Z} = \{\mathbf{X}, \mathbf{y}\}$  s.t. we obtain a real probability distribution:

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

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

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

#### 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 ??  $\theta^*$  that maximizes the posterior distribution.

#### 2.2.2. Maximization

**Definition 6.10**

**Maximum a-Posteriori Estimates (MAP):**

Is model/parameters  $\theta$  that maximize the posterior probability distribution:

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

**Log-MAP estimator:**

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

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

$$\stackrel{??}{\propto} \arg \max_{\theta} \{p(\mathbf{y}|\theta, \mathbf{X}) \cdot p_{\lambda}(\theta)\}$$

**Corollary 6.4 Negative Log MAP:**

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

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

### 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<sup>[def. 3.13]</sup>. 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<sup>[def. 3.10]</sup>, 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**  $p(w|y, X)/p(w|D)$   
**Posterior Probability Distribution:**  
 ① Specify the prior  $p_\lambda(w)$   
 ② Specify the likelihood  $p(y|w, X)/p(D|w)$   
 ③ Calculate the evidence  $p(y|X)/p(D)$   
 ④ Calculate the posterior distribution  $P(w|y, X)/p(w|D)$   

$$p(w|y, X) = \frac{p(y|w, X) \cdot p_\lambda(w)}{p(y|X)} = \frac{\text{Likelihood} \cdot \text{Prior}}{\text{Normalization}}$$

**Definition 6.13**  $p(y|X)/p(D)$   
**Marginal Likelihood** [see proof 10.2]: is the normalization constant that makes sure that the posterior distribution<sup>[def. 6.12]</sup> is an true probability distribution:  

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

**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<sup>[def. 6.12]</sup> of the observed data  $D = \{X, y\}$
- 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<sup>[def. 6.1]</sup> 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  $y = f(x) + \epsilon$ .
- Sometime 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

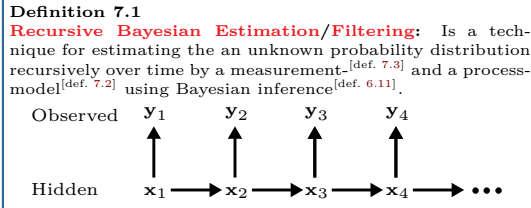


Figure 7: This problem corresponds to a *hidden Markov model* (HMM)<sup>[def. 14.1]</sup>  
 $x_t = (x_{t,1} \quad \dots \quad x_{t,n}) \quad y_t = (y_{t,1} \quad \dots \quad y_{t,m})$

**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**  $x_t \perp x_{1:t-2}|x_{t-1}$ : The process models<sup>[def. 7.2]</sup> is Markovian?? 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  $y_t$  to the model state  $x_t$

**Corollary 7.2**  $y_t \perp y_{1:t-1}x_{1:t-1}|x_t$   
**Conditional Independent Measurements:** The measurements  $y_t$  are conditionally independent of the previous observations  $y_{1:t-1}$  given the current state  $x_t$ :  

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

**Goal**  
 We want to combine the process model<sup>[def. 7.2]</sup> and the measurement model<sup>[def. 7.3]</sup> in a recursive way to obtain a good estimate of our model state:  

$$\left. \begin{array}{l} p(x_t|x_{t-1}) \\ p(y_t|x_t) \end{array} \right\} p(x_t|y_{1:t}) \xrightarrow{\text{recursion rule}} p(x_{t+1}|y_{1:t+1})$$

**Definition 7.4 Chapman-Kolmogorov eq.**  $p(x_t|y_{1:t-1})$   
**Prior Update/Prediction Step** [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** [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-filed do
3:   Prediction Step:

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

4:   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:

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

5: 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*?? process model<sup>[def. 7.2]</sup>,  $q$  with Gaussian model-noise  $v$  and a linear measurement model<sup>[def. 7.3]</sup>  $h$  with Gaussian process-noise  $w$ .

add difference of prediction, filtering and smoothing for posterior marginals

**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)

$$z[k] = H[k]x[k] + w[k-1] \quad \text{with} \quad w[k] \sim \mathcal{N}(0, R^{(k)})$$

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.

add Kalman algorithm (in slides Joseph Form 1 think)

# Gaussian Processes (GP)

## 1. Gaussian Process Regression

add complexity! Only due to inverse

### 1.1. Gaussian Linear Regression

Given

- Linear Model with Gaussian Noise:  

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2 \mathbf{I}) \quad (8.1)$$

$$\mathbf{y} = f(\mathbf{x}) + \epsilon$$
 $\Rightarrow$  Gaussian Likelihood:  $\mathbf{p}(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma_n^2 \mathbf{I})$
- Gaussian Prior:  $\mathbf{p}(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_p)$

Sought

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

**Definition 8.1**  $\mathbf{p}(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \mathcal{N}(\mu_w, \Sigma_w^{-1})$  proof 10.7:  
**Posterior Distribution**  

$$\mu_w = \frac{1}{\sigma_n^2} \mathbf{x}_w^\top \Sigma_w^{-1} \mathbf{X} \mathbf{y} \quad \Sigma_w = \frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}^\top + \Sigma_p^{-1}$$

Note

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

**Definition 8.2**  $\mathbf{p}(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mu_*, \Sigma_*)$  proof 10.8:  
**Posterior Predictive Distribution**  

$$\mu_* = \frac{1}{\sigma^2} \mathbf{x}_*^\top \Sigma_*^{-1} \mathbf{X} \mathbf{y} \quad \Sigma_* = \mathbf{x}_*^\top \Sigma_w^{-1} \mathbf{x}_*$$
 (8.2)

### 1.2. Kernelized Gaussian Linear Regression

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

$$\mu_*$$
 (8.4)

**Definition 8.4 Gaussian Process:**

## 2. Model Selection

### 2.1. Marginal Likelihood

# Approximate Inference

Problem

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

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

• Normalization constants:  

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

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

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  $\mathbf{p}$  as well as possible.

## 1. Variational Inference

**Definition 9.2 Bayes Variational Inference:**

Given an unnormalized (posterior) probability distribution:

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

BVI seeks an *approximate* probability distribution  $q_\lambda$ , that is parameterized by a *variational parameter*  $\lambda$  and approximates  $\mathbf{p}(\theta|\mathbf{y})$  well.

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

### 1.1. Laplace Approximation

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

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

the distribution is given by:

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

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

with

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

Note

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

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

### 1.2. Black Box Stochastic Variational Inference

The most common way of finding  $q_\lambda$  is by minimizing the KL-divergence<sup>[def. 3.8]</sup> between our approximate distribution  $q$  and our true posterior  $\mathbf{p}$ :

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

Note

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

**Definition 9.5** [proof 10.12]

**ELBO-Optimization Problem:**

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

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

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

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

**Attention:** Sometimes people write simply  $\mathbf{p}$  for the posterior and  $\cdot(\cdot)$  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.

### 1.3. 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:  

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

### 1.3.1. Maximizing The ELBO

**Definition 9.7 Gradient of the ELBO Loss:**

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

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

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

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

Problem

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

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

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:

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

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.

### 1.4. The Reparameterization Trick

**Principle 9.1** [proof 10.14]

**Reparameterization Trick:** Let  $\phi$  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  $\theta$  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 ??:

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

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

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

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

**Definition 9.8** [example 10.3]

**Reparameterized ELBO Gradient**<sup>[def. 9.7]</sup>:

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

$$\nabla_\lambda L(\lambda) = \nabla_\lambda \text{ELBO}(\lambda) \quad (9.13)$$

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

$$= \mathbb{E}_{\epsilon \sim \phi} [\nabla_\lambda \log \mathbf{p}(\mathbf{y}|g(\epsilon; \lambda))] - \nabla_\lambda \text{KL}(q_\lambda(\theta) \parallel \mathbf{p}(\theta))$$

**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  $\phi(\epsilon)$ :

$$\theta \sim q(\theta|\lambda) = \mathcal{N}(\theta; \mu, \Sigma) \quad \Rightarrow \quad \lambda = [\mu \quad \Sigma]$$

$$\epsilon \sim \phi(\epsilon) = \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$$

Then it follows for the ELBO:

$$\nabla_\lambda L(\lambda) = \nabla_\lambda \text{ELBO}(\lambda) \quad (9.14)$$

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

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

$$= \nabla_{\mathbf{C}, \mu} \text{KL}(q_{\mathbf{C}, \mu} \parallel \mathbf{p}(\theta))$$

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

$$= \nabla_{\mathbf{C}, \mu} \text{KL}(q_{\mathbf{C}, \mu} \parallel \mathbf{p}(\theta))$$

## 2. Markov Chain Monte Carlos Methods

**Definition 9.9**

**Markov Chain Monte Carlo (MCMC) Methods:**

## 3. Integrated Nested Laplace Approximation

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

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

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

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

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

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

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

Note

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

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

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

$$\Rightarrow \quad \tilde{\mathbf{p}}(\theta|\mathbf{y}) = \frac{\mathbf{p}(\mathbf{x}, \theta, \mathbf{y})}{\tilde{\mathbf{p}}(\mathbf{x}|\theta, \mathbf{y}) \mathbf{p}(\mathbf{y})} \propto \frac{\mathbf{p}(\mathbf{x}, \theta, \mathbf{y})}{\tilde{\mathbf{p}}_G(\mathbf{x}|\theta, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*(\theta)}$$

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

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

and then numerical integration w.r.t.  $\theta$ :

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

Note

$\tilde{\mathbf{p}}(\theta|\mathbf{y})$  is usually quiet 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):

- Model the prior over our weights  $\theta = [\mathbf{W}^0 \dots \mathbf{W}^L]$  by a neural network:

$$\theta \sim p_{\lambda}(\theta) = \mathbf{F} \quad \text{with} \quad \mathbf{F} = \mathbf{F}^L \circ \dots \circ \mathbf{F}^1$$

$$\mathbf{F}^l = \varphi \circ \mathbf{F}^l = \varphi(\mathbf{W}^l \mathbf{x} + \mathbf{b}^l)$$

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

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

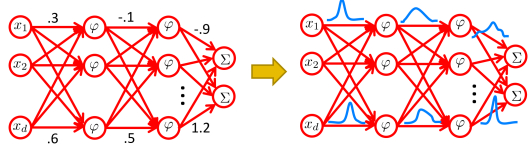


Figure 8

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

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

### Note

Recall for normal Bayesian Linear regression we had:

### Problem

All the weights of the prior  $p_{\lambda}(\theta) = \mathbf{F}$  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  $p(\theta|\mathbf{y}, \mathbf{X})$  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  $\mathbf{x}_i$  in order to obtain  $\mu(\mathbf{x}_i; \theta)$  and  $\sigma(\mathbf{x}_i; \theta)^2$ :

$$\theta^* = \arg \max_{\theta} \{p(\theta|\mathbf{X}, \mathbf{y})\} \stackrel{\text{eq. (6.13)}}{=} \arg \min_{\theta} \lambda \|\theta\|_2^2$$

$$- \sum_{i=1}^n \left( \frac{1}{2\sigma(\mathbf{x}_i; \theta)^2} \|y_i - \mu(\mathbf{x}_i; \theta)\|^2 + \frac{1}{2} \log \sigma(\mathbf{x}_i; \theta)^2 \right)$$

### Explanation 10.1. [def. 10.2]

- $\frac{1}{2} \log \sigma(\mathbf{x}_i; \theta)^2$ : tries to force neural network to predict small uncertainty
- $\frac{1}{2\sigma(\mathbf{x}_i; \theta)^2} \|y_i - \mu(\mathbf{x}_i; \theta)\|^2$ : 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 MAP Gradient of BNN: proof 10.15

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

### 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<sup>[def. 6.16]</sup>  $\iff$  it is overconfident in regions where we haven't even seen any data.  
 $\Rightarrow$  need to use Bayesian approach to approximate posterior distribution.

#### 0.1. Variational Inference For BNN

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

## 0.2. Making Predictions

### Proposition 10.1 Title:

## 1. Proofs

Proof 10.1: Definition 6.15:

$$p(\mathbf{f}_{*}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{f}_{*}, \mathbf{x}_{*}, \mathbf{X}, \mathbf{y})}{p(\mathbf{x}_{*}, \mathbf{X}, \mathbf{y})}$$

$$= \frac{\int p(\mathbf{f}_{*}, \mathbf{x}_{*}, \mathbf{X}, \mathbf{y}, \mathbf{w}) d\mathbf{w}}{\int p(\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}) d\mathbf{w}}$$

$$\stackrel{??}{=} \frac{\int p(\mathbf{f}_{*}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}, \mathbf{w}) p(\mathbf{w}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}) d\mathbf{w}}{\int p(\mathbf{f}_{*}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}, \mathbf{w}) p(\mathbf{w}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}) d\mathbf{w}}$$

$$\stackrel{??}{=} \frac{\int p(\mathbf{f}_{*}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}, \mathbf{w}) p(\mathbf{w}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}) d\mathbf{w}}{\int p(\mathbf{f}_{*}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}, \mathbf{w}) p(\mathbf{w}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}) d\mathbf{w}}$$

$$= \int p(\mathbf{f}_{*}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}, \mathbf{w}) p(\mathbf{w}|\mathbf{x}_{*}, \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

$$\stackrel{\clubsuit}{=} \int p(\mathbf{f}_{*}|\mathbf{x}_{*}, \mathbf{w}) p(\mathbf{w}|\mathbf{X}, \mathbf{y}) d\mathbf{w}$$

### Note ♣

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

Proof 10.2: Definition 6.13:

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}, \mathbf{w}|\mathbf{X}) d\mathbf{w} = \int p(\mathbf{y}|\mathbf{w}, \mathbf{X}) p(\mathbf{w}|\mathbf{X}) d\mathbf{w}$$

$$\stackrel{\text{eq. (6.6)}}{=} \int p(\mathbf{y}|\mathbf{w}, \mathbf{X}) p(\mathbf{w}) d\mathbf{w}$$

Proof 10.3: Definition 7.4:

$$p(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{y}_{1:t_1}) \stackrel{??}{=} p(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y}_{1:t_1}) p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t_1})$$

independ.  $p(\mathbf{x}_t|\mathbf{x}_{t-1}) p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t_1})$

marginalization/integration over  $\mathbf{x}_{t-1}$  gives the desired result.

Proof 10.4: Definition 7.5:

$$p(\mathbf{x}_t, \mathbf{y}_t|\mathbf{y}_{1:t-1}) \stackrel{??}{=} \begin{cases} p(\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{1:t-1}) p(\mathbf{y}_t|\mathbf{y}_{1:t-1}) \\ p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{y}_{1:t-1}) p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) \end{cases}$$

$$p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{y}_{1:t-1}) \stackrel{[\text{cor. 7.2}]}{=} p(\mathbf{y}_t|\mathbf{x}_t)$$

from which follows immediately eq. (7.5).

Proof 10.5: Definition 7.6:

$$p(\mathbf{y}_t|\mathbf{y}_{1:t-1}) = \int p(\mathbf{y}_t, \mathbf{x}_t|\mathbf{y}_{1:t-1}) d\mathbf{x}_t$$

$$= \int p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{y}_{1:t-1}) p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) d\mathbf{x}_t$$

$$\stackrel{[\text{cor. 7.2}]}{=} \int p(\mathbf{y}_t|\mathbf{x}_t) p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) d\mathbf{x}_t$$

Proof 10.6: [cor. 7.3]:

$$p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t}) \stackrel{??}{=} p(\mathbf{y}_{1:t}|\mathbf{x}_{1:t}) p(\mathbf{x}_{1:t})$$

$$\stackrel{??}{=} p(\mathbf{y}_{1:t}|\mathbf{x}_{1:t}) p(\mathbf{x}_t|\mathbf{x}_{t-1:t}) \dots p(\mathbf{x}_2|\mathbf{x}_1) p(\mathbf{x}_1)$$

$$\stackrel{\text{eq. (7.1)}}{=} p(\mathbf{y}_{1:t}|\mathbf{x}_{1:t}) \left( p(\mathbf{x}_1) \prod_{i=2}^t p(\mathbf{x}_i|\mathbf{x}_{i-1}) \right)$$

$$\stackrel{[\text{cor. 7.2}]}{=} p(\mathbf{y}_1|\mathbf{x}_1) \dots p(\mathbf{y}_t|\mathbf{x}_t) \left( p(\mathbf{x}_1) \prod_{i=2}^t p(\mathbf{x}_i|\mathbf{x}_{i-1}) \right)$$

$$= p(\mathbf{y}_1|\mathbf{x}_1) p(\mathbf{x}_1) \prod_{i=2}^t p(\mathbf{y}_i|\mathbf{x}_i) p(\mathbf{x}_i|\mathbf{x}_{i-1})$$

Proof 10.7: GP Posterior Distribution<sup>[def. 8.1]</sup>

$$p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w}) p(\mathbf{w})$$

$$\propto \exp \left( -\frac{1}{2} \frac{1}{\sigma_n^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \right) \exp \left( -\frac{1}{2} \mathbf{w}^T \Sigma^{-1} \mathbf{w} \right)$$

$$\propto \exp \left\{ -\frac{1}{2} \frac{1}{\sigma_n^2} \left( \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}^T \mathbf{w} + \sigma_n^2 \mathbf{w}^T \Sigma^{-1} \mathbf{w} \right) \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} \frac{1}{\sigma_n^2} \left( \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T (\mathbf{X}^T \mathbf{X}^T + \sigma_n^2 \Sigma^{-1}) \mathbf{w} \right) \right\}$$

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

$$p(\mathbf{w}|\mathcal{D}) \propto \exp \left( -\frac{1}{2} (\mathbf{w} - \mu_w)^T \Sigma_w (\mathbf{w} - \mu_w) \right)$$

$$\propto \exp \left( -\frac{1}{2} \left( \mathbf{w}^T \Sigma_w \mathbf{w} - 2\mathbf{w}^T \Sigma_w \mu_w + \mu_w^T \Sigma_w \mu_w \right) \right)$$

$\Sigma_w$  follows directly  $\Sigma_w = \sigma_n^{-2} \mathbf{X}^T \mathbf{X} + \Sigma_p$

$\mu_w$  follows from  $2\mathbf{w}^T \mathbf{X}^T \mathbf{y} = 2\mathbf{w}^T \Sigma_w \mu_w \Rightarrow \mu_w = \Sigma_w^{-1} \mathbf{X}^T \mathbf{y}$ .

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:

$$L_n(\theta) = \log(p(\theta|\mathbf{y})) = (\log \text{Prior} + \log \text{Likelihood})$$

we now approximate  $L_n(\theta)$  by a Taylor approximation around its maximum  $\hat{\theta}$ :

$$L_n(\theta) = L_n(\hat{\theta}) + \frac{1}{2} \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} (\theta - \hat{\theta})^2 + \mathcal{O}((\theta - \hat{\theta})^3)$$

we can no derive the distribution:

$$p(\theta|\mathbf{y}) \approx \exp(L_n(\theta)) = \exp(\log p(\theta|\mathbf{y}))$$

$$= p(\hat{\theta}) \exp \left( \frac{1}{2} \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} (\theta - \hat{\theta})^2 \right)$$

$$= \sqrt{2\pi\sigma^2} p(\hat{\theta}) \mathcal{N}(\theta; \hat{\theta}, \sigma) \approx \frac{1}{\sqrt{2\pi\sigma^2}} \mathcal{N}(\theta; \hat{\theta}, \sigma)$$

### 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:

$$\nabla L_n(\theta) = \nabla L_n(\theta_1, \theta_2) = 0$$

$$L_n(\theta) = L_n(\hat{\theta}) + \frac{1}{2} (A(\theta_1 - \hat{\theta}_1)^2 + B(\theta_2 - \hat{\theta}_2)^2 + C(\theta_1 - \hat{\theta}_1)(\theta_2 - \hat{\theta}_2))$$

$$L_n(\theta) = L_n(\hat{\theta}) + (\theta - \hat{\theta})^T H(\hat{\theta}) (\theta - \hat{\theta})$$

$$= L_n(\hat{\theta}) + \frac{1}{2} Q(\theta)$$

$$A = \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} \quad B = \frac{\partial^2 L_n}{\partial \theta^2} \Big|_{\hat{\theta}} \quad C = \frac{\partial^2 L_n}{\partial \theta_1 \partial \theta_2} \Big|_{\hat{\theta}}$$

$$H = \begin{bmatrix} A & C \\ C & B \end{bmatrix} \quad \Sigma = H^{-1}(\hat{\theta})$$

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

$$L_n(\theta) \approx L_n(\hat{\theta}) + (\theta - \hat{\theta})^T \nabla \nabla^T L_n(\hat{\theta}) (\theta - \hat{\theta})$$

$$H(\theta) = \nabla \nabla^T L_n(\theta) \quad \Sigma = H^{-1}(\hat{\theta})$$

$$p(\theta|\mathbf{y}) = \sqrt{(2\pi)^n \det(\Sigma)} p(\hat{\theta}) \mathcal{N}(\theta; \hat{\theta}, \Sigma)$$

$$\approx c \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \mathcal{N}(\theta; \hat{\theta}, \Sigma)$$

Proof 10.12: [def. 9.5]

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

Proof 10.13: [def. 9.6]

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

Proof 10.14: principle 9.1 Let:

$$\begin{aligned}
\epsilon \sim \phi(\epsilon) &\quad \text{correspond to} \quad X \sim f_X \\
\theta = g(\epsilon; \lambda) &\quad \mathcal{Y} = \{y|y = g(x), \forall x \in \mathcal{X}\}
\end{aligned}$$

then it follows immediately with ??:

$$\begin{aligned}
\theta \sim q_\lambda(\theta) &= q(\theta|\lambda) = \frac{f_X(g^{-1}(y))}{\left| \frac{dg}{dx}(g^{-1}(y)) \right|} \\
&= \phi(\epsilon) |\nabla_\epsilon g(\epsilon; \lambda)|^{-1} \\
\Rightarrow &\text{parameterized in terms of } \epsilon
\end{aligned}$$

Proof 10.15: [def. 10.3]

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

## 2. Examples

**Example 10.1 Laplace Approximation**  
**Logistic Regression Likelihood + Gaussian Prior:**

**Example 10.2 ELBO Bayesian Logistic Regression:**  
Suppose:

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

Then it follows for the terms of the ELBO:

$$\text{KL}(q_\lambda \parallel p(\theta)) = \frac{1}{2} \sum_{i=1}^d \left( \mu_i^2 + \sigma_i^2 - 1 - \ln \sigma_i^2 \right)$$

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

**Example 10.3 ELBO Gradient Gaussian:** Suppose:

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

we can reparameterize using principle 9.1 by using:

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

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

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

from this it follows:

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

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

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

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

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

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

**Example 10.4 BNN Likelihood Function Examples:**

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

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

## 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**  $\phi$ : 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<sup>[def. 11.10]</sup>, 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?? 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 they 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  $\mathcal{X} \subseteq \mathbb{R}^d$ .  
**Let  $S$**  be any finite subset of data  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathcal{X}$ .  
Then the kernel matrix  $\mathcal{K} \in \mathbb{R}^{n \times n}$  is defined by:  
$$\mathcal{K} = \phi(\mathbf{X})\phi(\mathbf{X}^\top) = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))(\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))^\top$$
$$= \begin{pmatrix} \mathbf{k}(\mathbf{x}_1, \mathbf{x}_1) & \dots & \mathbf{k}(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ \mathbf{k}(\mathbf{x}_n, \mathbf{x}_1) & \dots & \mathbf{k}(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix} = \begin{pmatrix} \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_n) \end{pmatrix}$$
$$\mathcal{K}_{ij} = \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

## Corollary 11.3

**Kernel Eigenvector Decomposition:**  
For any symmetric matrix (Gram matrix  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)_{i,j=1}^n$ ) there exists an eigenvector decomposition:  
$$\mathcal{K} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top \quad (11.4)$$

$\mathbf{V}$ : orthogonal matrix of eigenvectors  $(\mathbf{v}_{t,i})_{i=1}^n$   
 $\mathbf{\Lambda}$ : diagonal matrix of eigenvalues  $\lambda_i$

**Assuming** all eigenvalues  $\lambda_t$  are non-negative, we can calculate the mapping:

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

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

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

## 1.1. 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:**  $\mathbf{k}/\mathcal{K}$  must be symmetric:  
 $\mathbf{k}(\mathbf{x}, \mathbf{y}) = \mathbf{k}(\mathbf{y}, \mathbf{x}) = \phi(\mathbf{x})^\top \phi(\mathbf{y}) = \phi(\mathbf{y})^\top \phi(\mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$

**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 induce a *p.s.d. symmetric* kernel matrix  $\mathbf{k}$  for any possible  $S \subseteq \mathcal{X}$  see ?? 11.1.  
 $\iff$  all eigenvalues of the kernel gram matrix  $\mathcal{K}$  for *finite*  $\mathcal{V}$  must be non-negative ??.

## 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  $\phi$  s.t.

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

**Theorem 11.2 General Mercers Theorem:** Let  $\Omega$  be a compact subset of  $\mathbb{R}^n$ . Suppose  $\mathbf{k}: \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  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 condtions describe an inner product in a high dimensional space.  
 $\implies$  can replace the inner product by the kernel function.

Check if R or R\* as in script

## 3. The Kernel Trick

**Definition 11.6 Kernel Trick:** If a kernel has an analytic form we do no longer need to calculate:  
• the function mapping  $\mathbf{x} \mapsto \phi(\mathbf{x})$  and  
• the inner product  $\phi(\mathbf{x})^\top \phi(\mathbf{y})$   
explicitly but simply use the formula for the kernel:

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

see examples 11.1 and 11.2

## Note

- Possible to operate in any n-dimensional function space, efficiently.
- $\phi$  not necessary anymore.
- Complexity independent of the functions space.

## 4. Types of Kernels

### 4.1. Stationary Kernels

**Definition 11.7 Stationary Kernel:** A stationary kernel is a kernel that only considers vector differences:

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \mathbf{k}(\mathbf{x} - \mathbf{y}) \quad (11.10)$$

see example example 11.3

### 4.2. Isotropic Kernels

**Definition 11.8 Isotropic Kernel:** A isotropic kernel is a kernel that only considers distance differences:

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \mathbf{k}(\|\mathbf{x} - \mathbf{y}\|_2) \quad (11.11)$$

**Corollary 11.4 :**  
Isotropic  $\rightarrow$  Stationary

## 5. Important Kernels on $\mathbb{R}^d$

### 5.1. The Linear Kernel

**Definition 11.9 Linear/String Kernel:**  
$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y} \quad (11.12)$$

### 5.2. The Polynomial Kernel

**Definition 11.10 Polynomial Kernel:**  
represents all monomials?? of degree up to  $m$   
$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^\top \mathbf{y})^m \quad (11.13)$$

### 5.3. The Sigmoid Kernel

**Definition 11.11 Sigmoid/tanh Kernel:**  
$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \tanh \kappa \mathbf{x}^\top \mathbf{y} - b \quad (11.14)$$

### 5.4. The Exponential Kernel

**Definition 11.12 Exponential Kernel:**  
is an continuous kernel that is non-differential  $\mathbf{k} \in \mathcal{C}^0$ :  
$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|_1}{\theta}\right) \quad (11.15)$$

$\theta \in \mathbb{R}$ : corresponds to a threshold.

### 5.5. The Gaussian Kernel

**Definition 11.13 Gaussian/Squared Exp. Kernel/Radial Basis Functions (RBF):**  
Is an infinte dimensional smooth kernel  $\mathbf{k} \in \mathcal{C}^\infty$  with some usefull 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  $\theta$ ).  $2\theta \in \mathbb{R}$  corresponds to a threshold that determines how close input values need to be in order to be considered similar:

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

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 admitable functions vary much more.

## 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-(infinte)-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.

## 5.6. The Matern Kernel

When looking at actual data/sample paths the smoothness of the Gaussian kernel<sup>[def. 11.13]</sup> is often a too strong assumption that does not model reality the same holds true for the non-smoothness of the exponential kernel<sup>[def. 11.12]</sup>. 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  $\nu$ :

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

$\nu, \rho \in \mathbb{R}_+$   $\nu$ : Smoothness

$\mathcal{K}_\nu$  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.

### 6.1. 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:

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

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

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \alpha \mathbf{k}_1(\mathbf{x}, \mathbf{x}') \quad \alpha \in \mathbb{R}_+ \quad (11.20)$$

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

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

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

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

**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:**  $\mathcal{K} = \Phi^\top \Phi \in \mathbb{R}^{n \times n}$ .  
$$\mathbf{v}^\top \mathcal{K} \mathbf{v} = \mathbf{v}^\top \Phi^\top \Phi \mathbf{v} = (\Phi \mathbf{v})^\top \Phi \mathbf{v} = \|\Phi \mathbf{v}\|_2^2 \geq 0$$



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  $\phi$  as:

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

Operation Count:

- 2 · 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:

$$\left\langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \right\rangle = \underbrace{\left\langle \mathbf{x}_i, \mathbf{x}_j \right\rangle^2}_{:= \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j)} = \langle \{x_{i1}, x_{i2}\}, \{x_{i1}, x_{i2}\} \rangle^2$$
$$= (x_{i1} x_{i2} + x_{j1} x_{j2})^2$$
$$= x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2x_{i1} x_{i2} x_{j1} x_{j2}$$

Operation Count:

- 2 multiplicaitons 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*<sup>[def. 12.5]</sup> that influence the propagation of to the state variables  $\mathbf{x}$ .

**Definition 12.4 Output/Measurement Variables/State Observations:**  $\mathbf{y} \in \mathcal{O}$   
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 Measurement/Output/Observation Model**  $h$ :  
Describes the mapping of the state  $\mathbf{x}$  onto the output  $\mathbf{y}$ .

**Definition 12.7 (Discrete) State Space Model:**  
 $\mathbf{x}^{k+1} = f(t, \mathbf{x}^k, \mathbf{u}^k) \quad t = 1, \dots, K$  (12.2)  
 $\mathbf{y}^k = h(t, \mathbf{x}^k, \mathbf{u}^k)$  (12.3)

**Definition 13.6  $n^{\text{th}}$  Transition Probability**  $\mathbf{p}_{ji}^{(n)}(t)$ :  
denotes the probability of reaching state  $s_j$  from state  $s_i$  in  $n$  steps:  
 $\mathbf{p}_{ij}^{(n)}(t) = \mathbb{P}(X_{t+n} = s_j | X_t = s_i) \quad \forall s_i, s_j \in \mathcal{S} \quad (13.4)$

**Definition 13.7 Transition Matrix**  $\mathbf{P}(t)$ :  
The transition probabilities eq. (13.4) can be represented by a *row-stochastic matrix*??  $\mathbf{P}(t)$  where the  $i^{\text{th}}$  row represents the transition probabilities for the  $i^{\text{th}}$  state  $s_i$  i.e.

From  $s_i$  To  $s_j$

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:  
 $\sum \delta^+(s_i) = 1$  (13.5)

### 1.1. Simulating Markov Chains

**Corollary 13.2** proof 13.1  
**Realization of a Markov Chain:**  
 $\mathbf{P}(X_0 = x_0, \dots, X_N = x_N) = q_0(x_1) \prod_{n=1}^N \mathbf{p}_{n-1,n}(t)$

**Algorithm 13.1 Forward Sampling:**  
**Input:**  $\mathbf{q}(x_0)$  and  $\mathbf{P}$   
**Output:**  $\mathbf{P}(X_{0:N})$   
Sample  $x_0 \sim \mathbf{P}(X_0)$   
**for**  $j = 1, \dots, n$  **do**  
     $x_j \sim \mathbb{P}(X_j | X_{j-1} = x_{j-1})$   
**5: end for**

### 1.2. State Distributions

**Definition 13.8 Probability Distribution of the States**  $\mathbf{q}_{n+1}$ :  
 $q_{n+1}(s_j) = \mathbb{P}(X_{n+1} = s_j) \quad \forall s_i \in \mathcal{S}$   
 $= \sum_{i=1}^n \mathbb{P}(X_n = s_i) \mathbb{P}(X_{n+1} = s_j | X_n = s_i)$   
 $= \sum_{i=1}^n q_n(s_i) \mathbf{p}_{ij}(t)$  (13.6)  
 $\mathbf{q}_{n+1} = [q_{n+1}(s_1) \dots q_{n+1}(s_n)]$   
 $= \mathbf{q}_n \mathbf{P}(t)$   
 $= [q_n(s_1) \dots q_n(s_n)] \begin{bmatrix} \mathbf{p}_{1,1} & \mathbf{p}_{1,2} & \dots & \mathbf{p}_{1,n} \\ \mathbf{p}_{2,1} & \mathbf{p}_{2,2} & \dots & \mathbf{p}_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{p}_{n,1} & \mathbf{p}_{n,2} & \dots & \mathbf{p}_{n,n} \end{bmatrix} (t)$

**Corollary 13.3** [proof 13.2]  
**Time-homogeneous Markov Transition Probabilities:**  
 $\mathbf{q}_{n+1} = \mathbf{q}_0 \mathbf{P}^{n+1}$  (13.7)

**Definition 13.9 Stationary Distribution:**  
A markov chain has a stationary distribution if it satisfies:  
 $\lim_{N \rightarrow \infty} q_N(s_i) = \lim_{N \rightarrow \infty} \mathbb{P}(X_N = s_i) = \pi_i \quad \forall s_i \in \mathcal{S}$   
 $\lim_{N \rightarrow \infty} \mathbf{q}_N = [\pi_1 \dots \pi_n] \iff \mathbf{q} = \mathbf{q} \mathbf{P}(N)$  (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*!

edit matrix version with eigenvector  
add recurrent and transient states

### 1.3. Properties of States

**Definition 13.10 Absorbing State/Sink:** Is a state  $s_i$  that once entered cannot be left anymore:  
 $\mathbf{p}_{ij}^{(n)}(t) = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$  (13.9)

**Definition 13.11 Accessible State**  $s_i \rightarrow s_j$ :  
A state  $s_j$  is accessible from state  $s_i$  iff:  
 $\exists n : \mathbf{p}_{ij}^{(n)}(t) > 0$  (13.10)

**Definition 13.12 Communicating States**  $s_i \leftrightarrow s_j$ :  
Two states  $s_j$  and  $s_i$  are communicating iff:  
 $\exists n_1 : \mathbf{p}_{ij}^{(n_1)}(t) > 0 \quad \wedge \quad \exists n_2 : \mathbf{p}_{ji}^{(n_2)}(t) > 0$  (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:  
 $k = \gcd\{n > 0 : \mathbf{p}_{ii}^{(n)} = \mathbb{P}(X_n = s_i | X_0 = s_i) > 0\}$  (13.12)

**Definition 13.14 Aperiodic State**  $k = 1$ :  
Is a state  $s_i$  with periodicity<sup>[def. 13.13]</sup> of one  $\iff k = 1$

**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.*

### 1.4. Characteristics of Markov Processes/Chains

**Definition 13.15 Time-homogeneous/Stationary Markov Chain:**  
are markov chains<sup>[def. 13.3]</sup> where the transition probability is independent of time:  
 $\mathbf{p}_{ji} = \mathbb{P}(X_t = s_j | X_{t-1} = s_i) = \mathbb{P}(X_{t-\tau} = s_j | X_{t-\tau} = s_i) \quad \forall \tau \in \mathbb{N}_0$  (13.13)

**Corollary 13.7**  $\mathbf{P}$   
**Transition Matrices of Stationary MCs:**  
Transition matrices of time-homogeneous markov chain are constant/time independent:  
 $\mathbf{P}(t) = \mathbf{P}$  (13.14)

**Definition 13.16 Aperiodic Markov Chain:** Is a markov chain where all states are aperiodic:  
 $\gcd\{n > 0 : \mathbf{p}_{ii}^{(n)} = \mathbb{P}(X_n = s_i | X_0 = s_i) > 0\} = 1$   
 $\forall i \in \{1, \dots, n\}$  (13.15)

**Definition 13.17 Irreducible Markov Chain:** Is a Markov chain that has only *communicating states*<sup>[def. 13.12]</sup>:  
 $s_j \leftrightarrow s_i \quad \forall i, j \in \{1, \dots, n\}$  (13.16)  
•  $\implies$  no sinks<sup>[def. 13.10]</sup>  
•  $\implies$  every state can be reached from every other state

**Corollary 13.8** : An *irreducible*<sup>[def. 13.17]</sup> markov chain is automatically *aperiodic*<sup>[def. 13.16]</sup> if it has at least one aperiodic state<sup>[def. 13.14]</sup>  $\iff$  *ergodic*<sup>[def. 13.18]</sup>.

**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<sup>[def. 13.17]</sup>
- 2 Aperiodic<sup>[def. 13.16]</sup>

**Corollary 13.10 Stationary Distribution:** An ergodic markov chain has a *unique* stationary distribution<sup>[def. 13.9]</sup> and converges to it starting from any initial state  $q_0(s_i)$

### 1.5. Types of Markov Chains

	Observable	Unobservable	
Uncontrolled	MC <sup>[def. 13.3]</sup>	HMM <sup>[def. 14.1]</sup>	
Controlled	MDP <sup>[def. 15.1]</sup>	POMDP <sup>[def. 16.1]</sup>	

### 1.6. Markov Chain Monte Carlo (MCMC)

### 2. Proofs

Proof 13.1: <sup>[cor. 13.2]</sup>  
 $\mathbb{P}(X_0 = x_0, \dots, X_N = x_N) = \mathbb{P}(X_0 = x_0) \cdot \mathbb{P}(X_1 = x_1 | X_0 = x_0) \cdot \mathbb{P}(X_2 = x_2 | X_1 = x_1, X_0 = x_0) \cdot \dots \cdot \mathbb{P}(X_N = x_N | X_{N-1} = x_{N-1}, \dots, X_0 = x_0)$   
and then simply use the Markovian property

Proof 13.2: Corollary 13.3

$$\mathbf{q}_{n+1} = \mathbf{P} \mathbf{q}_n = (\mathbf{q}_{n-1} \mathbf{P}) \mathbf{P} = \mathbf{q}_0 \mathbf{P}^{n+1}$$

### 3. Examples

**Example 13.1 Ergodic Markov Chain:**

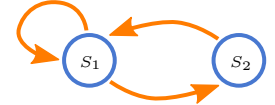
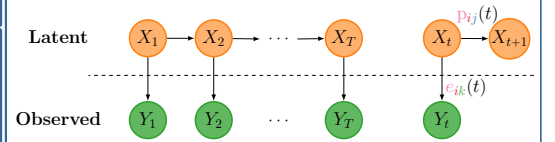


Figure 9: Ergodic for  $N = 2$  (can reach  $s_2$  at any  $t \leq N$  after  $N = 2$ )

## Hidden Markov Model (HMM)

**Definition 14.1**  $(\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathbf{P}, \mathbf{E})$   
**Hidden Markov Model (HMM):**  
Is a Markov Chain<sup>[def. 13.3]</sup> with hidden/latent states  $S_j$  that are only partially observable by noisy/indirect observations<sup>[def. 14.2]</sup>. It is characterized by the 5-tuple of:  
① States<sup>[def. 13.1]</sup>  $\mathcal{S} = \{s_1, \dots, s_n\}$   
② Actions<sup>[def. 15.2]</sup>  $\mathcal{A}/\mathcal{A}_{s_j} = \{a_1, \dots, a_m\}$   
③ Observations<sup>[def. 14.2]</sup>  $\mathcal{O}/\mathcal{O}_{s_j} = \{o_1, \dots, o_m\}$   
④ Transition Probabilities<sup>[def. 13.5]</sup>  $\mathbf{P}(s_i, s_j)$   
⑤ Emission/Output Probabilities<sup>[def. 14.3]</sup>  $e_{ij}(t)$



**Definition 14.2 Observations**  $\mathcal{O} = \{o_1, \dots, o_l\}$ :  
Are indirect or noisy observations that are related to the true states  $s_j$ .

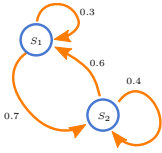
**Definition 14.3 Emission/ Output Probabilities**  $e_{ij}(t)$ :  
Given a state  $X_t = s_i$  the output probability is the probability of the output random variable  $Y_t$  to be in state  $o_j$ :

$$e_{ij}(t) = \mathbb{P}(Y_t = o_j | X_t = s_i) \quad \begin{cases} \forall o_i \in \mathcal{O} \\ \forall s_j \in \mathcal{S} \end{cases} \quad (14.1)$$

### 1. Markov Chains

**Definition 13.3 Markov Chain:**  
Is a sequence of random variables  $\{X_i\}_{i \in \mathcal{T}}$ ?? that processes the markovian property<sup>[def. 13.2]</sup> 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**  $\mathbf{q}_0$ : Describes the initial distribution of states:  
 $q_0(s_i) = \mathbb{P}(X_0 = s_i) \quad \forall s_i \in \mathcal{S}$   
 $\iff \mathbf{q}_0 = [q_0(s_1) \dots q_0(s_n)]$  (13.2)

**Definition 13.5 Transition Probability**  $\mathbf{p}_{ji}(t)$ :  
is the probability of a random variable  $X_t$  in state  $s_i$  to transition into state  $s_j$ :  
 $\mathbf{p}_{ij}(t) = \mathbb{P}(X_{t+1} = s_j | X_t = s_i) \quad \forall s_i, s_j \in \mathcal{S}$  (13.3)

# Markov Decision Processes (MDP)

**Definition 15.1**  $(S, \mathcal{A}, P_a, R_a)$   
**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<sup>[def. 13.1]</sup>  $S = \{s_1, \dots, s_n\}$
- ② Actions<sup>[def. 15.2]</sup>  $\mathcal{A}/\mathcal{A}_{s_j} = \{a_1, \dots, a_m\}$
- ③ Transition Probabilities<sup>[def. 15.3]</sup>  $P_a(s_i, s_j)$
- ④ Rewards<sup>[def. 15.4]</sup>  $r_a(s_i, s_j)$

**Definition 15.2**  
**Actions**  $\mathcal{A}_{s_i} = \{a_1, \dots, a_m\}$ :  
 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**  $P_a(s_j, s_i)(t)$ :  
 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  $a$ :

$$P_a(s_j, s_i) = P(s_j | s_i, a) = P(X_{t+1} = s_j | X_t = s_i, a_t = a) \quad \forall s_i, s_j \in S, \forall a \in \mathcal{A} \quad (15.1)$$

**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<sup>[def. 15.4]</sup>, 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<sup>[def. 15.8]</sup> and can be deterministic or randomized:

$$\pi : S \mapsto \mathcal{A} \quad \text{or} \quad \pi : S \mapsto \mathbb{P}(\mathcal{A}) \quad (15.3)$$

**Definition 15.6 Discounting Factor**  $\gamma$ :  
 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<sup>[def. 15.4]</sup> of the whole markov process given an initial state  $X_0 = x$ :

$$V^{\pi}(x) = J(\pi) | X_0 = x = \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \middle| X_0 = x \right] \quad (15.5)$$

$$= \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \middle| X_0 = x \right] \quad (15.6)$$

$$(15.7)$$

## 1.1. Calculating the value of $V^{\pi}$

**Definition 15.9** [proof 16.1]  
**Value Iteration:**  

$$V^{\pi}(x) = J(\pi) | X_0 = x \quad (15.8)$$

$$= \mathbb{E}_{x'|x, \pi(x)} [r(x, \pi(x)) + \gamma V^{\pi}(x')] = r(x, \pi(x)) + \gamma \mathbb{E}_{x'|x, \pi(x)} [V^{\pi}(x')] = r(x, \pi(x)) + \gamma \sum_{x' \in S} \mathbb{P}(x' | x, \pi(x)) V^{\pi}(x')$$

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:

$$V^{\pi} = \begin{bmatrix} V^{\pi}(s_1) \\ \vdots \\ V^{\pi}(s_n) \end{bmatrix} \quad \mathbf{r}^{\pi} = \begin{bmatrix} r^{\pi}(s_1, \pi(s_1)) \\ \vdots \\ r^{\pi}(s_n, \pi(s_n)) \end{bmatrix}$$

$$P^{\pi} = \begin{bmatrix} \mathbb{P}(s_1 | s_1, \pi(s_1)) & \mathbb{P}(s_2 | s_1, \pi(s_1)) & \dots & \mathbb{P}(s_n | s_1, \pi(s_1)) \\ \mathbb{P}(s_1 | s_2, \pi(s_2)) & \mathbb{P}(s_2 | s_2, \pi(s_2)) & \dots & \mathbb{P}(s_n | s_2, \pi(s_2)) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{P}(s_1 | s_n, \pi(s_n)) & \mathbb{P}(s_2 | s_n, \pi(s_n)) & \dots & \mathbb{P}(s_n | s_n, \pi(s_n)) \end{bmatrix}$$

### 1.1.1. Direct Mehtods

**Corollary 15.1 LU-decomposition**  $\mathcal{O}(n^3)$ :  
 The linear system from eq. (15.9):  

$$(\mathbf{I} - \gamma \mathbf{P}^{\pi}) V^{\pi} = \mathbf{r}^{\pi} \quad (15.10)$$
 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( $\mathbf{P}^{\pi}$ )  $\leq 1$ .

### 1.1.2. Fixed Point Iteration

**Corollary 15.2 Fixed-Point Iteration**  $\mathcal{O}(n \cdot |S|)$ :  
 The linear system from eq. (15.9) can be solve using *fixed-point iteration*?? in at most  $\mathcal{O}(n \cdot |S|)$  (if every state  $s_i$  is connected to every other state  $s_j \in S$ )

**Algorithm 15.1 Fixed Point Iteration:**  
**Input:** Initial Guess:  $V_0^{\pi} \stackrel{\text{i.e.}}{=} 0$   
 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?? that leads to a fixed point  $V^{\pi}$  with a rate depending on the discount factor  $\gamma$ .  

$$\|V_t^{\pi} - V^{\pi}\| = \|\phi V_{t-1}^{\pi} - \phi V^{\pi}\| \leq \gamma \|V_{t-1}^{\pi} - V^{\pi}\| = \gamma^t \|V_0^{\pi} - V^{\pi}\| \quad (15.12)$$

**Explanation 15.2.**  

- $\gamma \downarrow$ : the less we plan ahead/the smaller we choose  $\gamma$  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  $\gamma$  the longer it takes to converge but we will explore all possibilities. But for to large  $\gamma$  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

## 1.2. Choosing The Policy

**Question** how should we choose the  $\pi$ ? **Idea** compute  $J(\pi)$  for every possible policy:  

$$\pi^* = \arg \max J(\pi) \quad (15.13)$$

**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  $J/V^{\pi}$  depend on  $\pi$  but if we do not know  $\pi$  yet we cannot compute those.

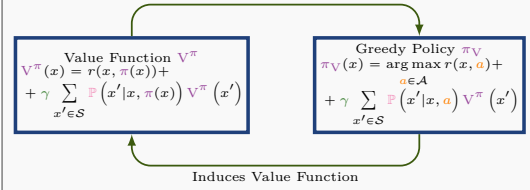
### 1.2.1. Greedy Policy

**Definition 15.10 Greedy Policy:**  
**Assuming** we know  $V^{\pi^{t-1}}$  then we could choose a greedy policy:  

$$a^* = \pi_t(x) \quad (15.14)$$

$$:= \arg \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in S} \mathbb{P}(x' | x, a) V^{\pi^{t-1}}(x')$$

- ① Given a policy  $\pi$  however we can calculate a value function  $V^{\pi}$
- ② Given a value function  $V$  we can induce a greedy policy<sup>[def. 15.10]</sup>  $\pi$  w.r.t.  $V$  Induces Policy



**Theorem 15.1 Optimality of Policies** [Bellman]:  
 A policy  $\pi_V$  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):

$$V^*(x) = \max_{a \in \mathcal{A}} \left[ r(x, a) + \gamma \sum_{x' \in S} \mathbb{P}(x' | x, a) V^*(x') \right] \quad (15.15)$$

$$:= \max_{a \in \mathcal{A}} Q^*(x, a) \quad (15.16)$$

### 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:  $\pi$   
 1: **while** Not converged  $t = t + 1$  **do**  
 2:   Compute  $V^{\pi^t}(x)$   

$$V^{\pi^t}(x) = r(x, \pi(x)) + \gamma \sum_{x' \in S} \mathbb{P}(x' | x, \pi_t(x)) V^{\pi^t}(x')$$
  
 3:   Compute greedy policy  $\pi_G$ :  

$$\pi_G(x) = \arg \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in S} \mathbb{P}(x' | x, a) V^{\pi^t}(x')$$
  
 4:   Set  $\pi_{t+1} \leftarrow \pi_G$   
 5: **end while**

### Algorithm 15.2

### Pros

- Monotonically improves  $V^{\pi^t} \geq V^{\pi^{t-1}}$
- is guaranteed to converge to an optimal policy/solution  $\pi^*$  in polynomial #iterations:  $\mathcal{O}\left(\frac{n^2 m}{1-\gamma}\right)$

### Cons

- Complexity *per iteration* requires to evaluate the policy  $V^{\pi}$  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:  

$$Q_t(x, a) = r(x, a) + \gamma \sum_{x' \in S} \mathbb{P}(x' | x, a) V_{t-1}(x') \quad \forall a \in \mathcal{A} \quad \forall x \in S$$
  
 3:   for all  $x \in S$  let:  

$$V_t(x) = \max_{a \in \mathcal{A}} Q_t(x, a)$$
  
 4:   **if**  $\max_{x \in S} |V_t(x) - V_{t-1}(x)| \leq \epsilon$  **then**  
 5:     break  
 6:   **end if**  
 7: **end for**  
 9: Choose greedy policy  $\pi_{V_t}$  w.r.t.  $V_t$

**Corollary 15.4** [proof 16.4]

**Value Iteration Contraction:**  
 Algorithm 15.3 is guaranteed to converge to a  $\epsilon$  optimal policy:

$$\left\| (V_t - V^*) \right\|_{\infty} \leq \gamma^t \left\| (V_0 - V^*) \right\|_{\infty} \quad (15.17)$$

$$\implies t \approx \ln \frac{\gamma}{\epsilon} \left\| V_0 - V^* \right\|_{\infty} \quad \text{for} \quad \left\| (V_t - V^*) \right\|_{\infty} \leq \epsilon$$

### Algorithm 15.3

### Pros

- Finds  $\epsilon$ -optimal solution in polynomial #iterations  $\mathcal{O}(\ln \frac{1}{\epsilon})$ <sup>[cor. 15.4]</sup>.
- Complexity *per iteration* requires us to solve a linear system  $\mathcal{O}(m \cdot n \cdot s) = \mathcal{O}(|\mathcal{A}| \cdot |S| \cdot s)$  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  $\mathcal{O}(n) = \mathcal{O}(|S|)$

### Cons

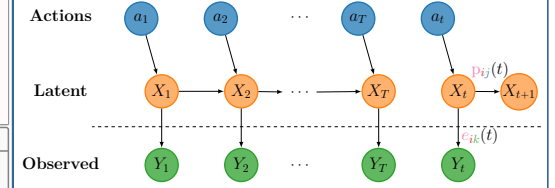
- Only  $\epsilon$ -optimal solution.

## Partially Observable MDP (POMDP)

**Definition 16.1**  $(S, \mathcal{A}, \mathcal{O}, P_a, E, R_a)$

**Partially Observable Markov Decision Process:**  
 A (POMDP) is a markov decision process<sup>[def. 15.1]</sup> with hidden markov states<sup>[def. 14.1]</sup>. It is characterized by the 6-tuple of:

- ① States<sup>[def. 13.1]</sup>  $S = \{s_1, \dots, s_n\}$
- ② Actions<sup>[def. 15.2]</sup>  $\mathcal{A}/\mathcal{A}_{s_j} = \{a_1, \dots, a_m\}$
- ③ Observations<sup>[def. 14.2]</sup>  $\mathcal{O}/\mathcal{O}_{s_j} = \{o_1, \dots, o_m\}$
- ④ Transition Probabilities<sup>[def. 15.3]</sup>  $P_a(s_i, s_j)$
- ⑤ Emission/Output Probabilities<sup>[def. 14.3]</sup>  $e_{ij}(t)$
- ⑥ Rewards<sup>[def. 15.4]</sup>  $r_a(s_i, s_j)$



### 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<sup>[def. 15.1]</sup> by introducing a *belief state space*  $\mathcal{B}$ .

**Definition 16.2 History**  $H_t$ :

Is a sequence of actions, observations and rewards:  
 $H_t = \{\{a_0, o_0, r_0\}, \dots, \{a_{t-1}, o_{t-1}, r_{t-1}\}\}$

**Definition 16.3 Belief State Space**  $\mathcal{B}$ : Is a  $|\mathcal{S}| - 1$  dimensional simplex or  $(|\mathcal{S}| - 1)$ -dimensional probability vector?? 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  $H_t$ <sup>[def. 16.2]</sup>.

### 1.1. Transition Model

**Definition 16.5** [proof 16.5]

**POMDP State/Posterior Update:**  
 $b_{t+1}(s_i) = \mathbb{P}(X_{t+1} = s_i | Y_{t+1} = o_k)$

$$= \frac{1}{Z} \mathbb{P}(Y_{t+1} = o_k | X_{t+1} = s_i, a_t) \cdot \sum_{s_j \in \mathcal{S}} b_t(s_j) \mathbb{P}(X_{t+1} = s_i | X_t = s_j, a_t) \quad (16.2)$$

**Definition 16.6 Stochastic Observation Model:**

$$\mathbb{P}(Y_{t+1} = o_k | b_t, a_t) = \sum_{s_i \in \mathcal{S}} b_t(s_i) \mathbb{P}(Y_{t+1} = o_k | X_t = s_i, a_t) \quad (16.3)$$

### 1.2. Reward Function

**Definition 16.7 POMDP Reward Function:**

$$r(b_t, a_t) = \sum_{s_j \in \mathcal{S}} b_t(s_j) r(s_j, a_t) \quad (16.4)$$

#### Note

For finite horizon  $T$ , the set of reachable belief states is finite however exponential in  $T$ .

add definition of simplex to math appendix which is basically this.

## 2. Proofs

### 2.1. Markov Decision Processes

Proof 16.1: [def. 15.8]

$$\begin{aligned} V^\pi(x) &= \mathbb{E}_{X_{1:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \\ &= \mathbb{E}_{\mathbf{X}} \left[ \gamma^0 r(X_0, \pi(X_0)) + \sum_{t=1}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \end{aligned}$$

$$\begin{aligned} \stackrel{\gamma^0=1}{X_0=x} &= r(x, \pi(x)) + \mathbb{E}_{\mathbf{X}} \left[ \sum_{t=1}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \\ &\stackrel{\text{re-index}}{=} r(x, \pi(x)) + \mathbb{E}_{\mathbf{X}} \left[ \sum_{t=0}^{\infty} \gamma^{t+1} r(X_{t+1}, \pi(X_{t+1})) \mid X_0 = x \right] \end{aligned}$$

$$\begin{aligned} &= r(x, \pi(x)) + \gamma \mathbb{E}_{\mathbf{X}} \left[ \sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_0 = x \right] \\ &= r(x, \pi(x)) \end{aligned}$$

$$+ \gamma \mathbb{E}_{X_1} \left[ \mathbb{E}_{X_{2:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_1 = x' \right] \mid X_0 = x \right]$$

$$\stackrel{??}{=} r(x, \pi(x))$$

$$+ \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, \pi(x)) \mathbb{E}_{X_{2:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_1 = x' \right]$$

$$\stackrel{\text{eq. (13.13)}}{=} r(x, \pi(x))$$

$$+ \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, \pi(x)) \mathbb{E}_{X_{2:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x' \right]$$

$$= r(x, \pi(x)) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, \pi(x)) \underline{V^\pi(x')}$$

Proof 16.2 [cor. 15.3]: Consider  $V, V' \in \mathbb{R}^n$  and let  $\phi$ :  
 $\phi x := r^\pi + \gamma P^\pi x \implies \phi V^\pi = V^\pi$

then it follows:

$$\begin{aligned} \left\| \phi V - \phi V' \right\| &= \left\| \cancel{\phi V} + \gamma P^\pi V - \cancel{\phi V'} - \gamma P^\pi V' \right\| \\ &= \left\| \gamma P^\pi (V - V') \right\| \\ &\stackrel{??}{\leq} \gamma \left\| P^\pi \right\| \cdot \left\| (V - V') \right\| \\ &\stackrel{\text{i.e. } L_2}{\leq} \gamma \cdot 1 \cdot \left\| (V - V') \right\|_2 \end{aligned}$$

Proof 16.3: algorithm 15.3

$$V_0(x) = \max_{a \in \mathcal{A}} r(x, a)$$

$$V_1(x) = \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, a) V_0(x')$$

$$V_{t+1}(x) = \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{S}} \mathbb{P}(x' | x, a) V_t(x')$$

Proof 16.4: [cor. 15.4] Let  $\phi : \mathbb{R}^n \mapsto \mathbb{R}^n$ , with:

$$(\phi V^*) (x) = Q(x, a) = \max_a \left[ r(x, a) + \gamma \sum_{x'} \mathbb{P}(x' | x, a) \right]$$

Bellman's theorem 15.1

and consider  $V, V' \in \mathbb{R}^n$

$$\phi V^* = V^*$$

$$\begin{aligned} \left\| \phi V - \phi V' \right\|_\infty &= \max_x \left| (\phi V)(x) - (\phi V')(x) \right| \\ &= \max_x \left| \max_a Q(x, a) - \max_{a'} Q'(x, a') \right| \\ &\stackrel{??}{\leq} \max_x \max_a \left| Q(x, a) - Q'(x, a) \right| \\ &= \max_{x, a} \left| f + \gamma \sum_{x'} \mathbb{P}(x' | x, a) V(x') - f' - \gamma \sum_{x'} \mathbb{P}(x' | x, a) V'(x') \right| \\ &= \gamma \max_{x, a} \left| \sum_{x'} \mathbb{P}(x' | x, a) (V(x') - V'(x')) \right| \\ &\stackrel{??}{\leq} \gamma \max_{x, a} \left| \sum_{x'} \mathbb{P}(x' | x, a) \right| \cdot \left| (V(x') - V'(x')) \right| \\ &\leq \gamma \cdot 1 \cdot \left\| (V(x') - V'(x')) \right\|_\infty \end{aligned}$$

#### Note

For the policy iteration the calculation was easier as the rewards canceled, however here we have the max.

### 2.2. MDPs

Proof 16.5: Defintion 16.5 Directly by definition 7.5 and its corresponding proof 10.4 with additional action  $a_t$ :

$$\begin{aligned} b_{t+1}(s_i) &= \mathbb{P}(X_{t+1} = s_i | y_{t+1}) \\ &= \frac{1}{Z} \mathbb{P}(y_{1:t+1} | s_i) \sum_{j=1} \frac{\mathbb{P}(X_{t+1} = s_i | y_{1:t})}{\underbrace{\mathbb{P}(X_t = s_j | y_{1:t})}_{b_t(s_j)} \mathbb{P}(s_i | s_j)} \end{aligned}$$

# Reinforcement Learning

Now we are working with an *unknown* MDP<sup>[def. 15.1]</sup> meaning that:

- ① we do no longer know the transition model<sup>[def. 15.3]</sup>
- 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**  $\tau$ :  
Is a set of consecutive 3-tuples of states, actions and rewards:  
 $\tau = \{s_t, a_t, r_t\} \quad t = 1, \dots, \tau \quad (17.1)$

**Definition 17.6 Episodic Learning:** Is a setting where we generate multiple  $K$ -episodes of different trajectories  $\{\tau^{(k)}\}_{i=1}^K$  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<sup>[def. 15.1]</sup> by:

- ① Estimating
  - the transition probabilities<sup>[def. 15.3]</sup>  $p_a(s_i, s_j)$
  - the reward function<sup>[def. 15.4]</sup>  $r(b_t, a_t)$
- ② Optimizing the policy of the estimated MDP

### 1.1. Estimating Transitions and Rewards

**Formula 17.1 Estimating Transitions and Rewards:**  
Given a data set  $D = \{(\mathbf{x}_0, a_0, r_0, \mathbf{x}_1), (\mathbf{x}_1, a_1, r_1, \mathbf{x}_2), \dots\}$  we estimate the transitions and rewards using a categorical distribution??:

$$N_{s_i|s_j,a} := \sum_{k=1}^t \delta(X_{k+1}=s_i|X_k=s_j, A_k=a) \quad (17.2)$$

$$N_{s_j,a} := \sum_{k=1}^t \delta(X_k=s_j, A_k=a) \quad (17.3)$$

$$p_a(s_i, s_j) \approx \frac{N_{s_i|s_j,a}}{N_{s_j,a}} \quad (17.4)$$

$$r(s_i, a) \approx \frac{1}{N_{s_i,a}} \sum_{k=1}^t \delta(X_k=s_i, A_k=a) r(X_k, A_k) \quad (17.5)$$

### 1.2. Choosing the next step

How should we choose the action  $a \in \mathcal{A}$  in order to balance exploration vs exploitation?

### 1.3. $\epsilon_t$ 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  $\epsilon_t$  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.

add general definition of RM conditions and sequence

**Pros**

- Simple
- Clearly sub optimal actions are not eliminated fast enough

**Cons**

### 1.4. The $R_{\max}$ Algorithm

**Algorithm 17.2** [Brafman & Tennenholz '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  $\Delta$  of observations for policy update

**Compute Policy**  $\pi_1$  of the MDP<sup>[def. 15.1]</sup> 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:
      
$$N_{\mathbf{x}_t, a} + = 1 \quad r(x_t, a) + = r(x_t, a) \quad (17.8)$$

      
$$N_{\mathbf{x}_{t+1}|\mathbf{x}_t, a} + = 1 \quad (17.9)$$

```

```
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**

**Cons**

- Memory: for all  $a \in \mathcal{A}$ ,  $\mathbf{x}_{t+1}, \mathbf{x}_t \in \mathcal{X}$  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  $\pi_t$  using policy (?? 1.2.2) or value iteration (?? 1.2.3)  $|\mathcal{A}| \cdot |\mathcal{S}|$  whenever we update out policy.

#### 1.4.1. How many transitions do we need?

**Proposition 17.2** [proof 17.1]  
**Number of Samples to bound Reward:**

$$\mathbb{P}(\hat{r}(s, a) - r(s, a) \leq \epsilon) \geq 1 - \delta \iff n \in \mathcal{O}\left(\frac{R_{\max}^2}{\epsilon^2} \log \frac{1}{\delta}\right) \quad (17.10)$$

**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  $\epsilon$ -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<sup>[def. 15.8]</sup> directly in order to act greedily upon it.

- Policy Gradient Methods
- Actor Critic Methods

### 2.1. Temporal Difference Learning (TD)

**Assume** we fix a random initial policy  $\pi$  and s.t. we have  $\hat{V}_0^{\pi}(s_j)$ .

**Goal:** want to calculate an unknown value function  $V^{\pi}$ .

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 carlos sampling?? 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  $\alpha_t$  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:

$$\mathbb{P}(\hat{V} \rightarrow \hat{V}^{\pi}) = 1 \quad (17.14)$$

### 2.2. 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

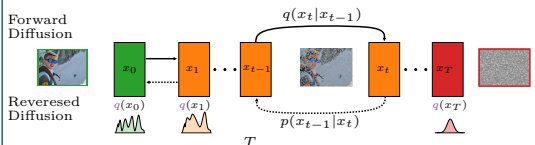
Proof 17.1: proposition 17.2 using hoeffdings bound?? with  $\delta$  and  $b - a = R_{\max}$ .

Diffusion Models

**Definition 17.9 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:

Forward Diffusion

Reversed Diffusion


$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \tag{17.16}$$

$q(\mathbf{x}_0)$ : true distribution of our input data:

Histroy

**Definition 17.10**  
**Forward Diffusion Process:**  
The forward diffusion process incrementally adds noise to the input:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}\left(\sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}\right) \quad \{\beta_t\}_{t=1}^T \in (0,1) \\ \beta_1 < \beta_2 < \dots < \beta_T$$
$$\mathbf{x}_t = \sqrt{1-\beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\epsilon \quad \epsilon \sim \mathcal{N}(0,1) \tag{17.17}$$

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.

$$\lim_{T \rightarrow \infty} q(\mathbf{x}_{1:T}|\mathbf{x}_0) \approx \mathcal{N}(0,\mathbf{I}) \tag{17.18}$$

**One Step Forward Process:**

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I}\right) \quad \bar{\alpha}_t := \prod_{s=0}^t \alpha_s$$
$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + (1-\bar{\alpha}_t)\epsilon \quad \epsilon \sim \mathcal{N}(0,\mathbf{I})$$

**Explanation 17.2.**  
*One Step Forward Diffusion Step :*  
*Sampling from a Gaussian and applying eq. (17.17) repeatedly to obtain  $q(\mathbf{x}_t|\mathbf{x}_0)$  using eq. (17.16) 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. (17.16).*

Notes

If the step-sizes  $\beta$  are too large it becomes to 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.

**Definition 17.11 Reverses Diffusion Process:**

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)) \tag{17.20}$$

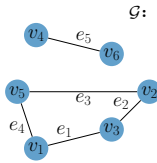
Latent Diffusion Models



# Graph Theory

## Definition 18.1 Graph

A graph  $\mathcal{G}$  is a pair  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  of a finite set of vertices  $\mathcal{V}$ <sup>[def. 18.4]</sup> and a multi set?? of edges  $\mathcal{E}$ <sup>[def. 18.8]</sup>.



## 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  $\mathcal{G}$ <sup>[def. 18.1]</sup> of order  $n$ .

**Corollary 18.2  $(p, q)$ -Graph:** Is a graph  $\mathcal{G}$ <sup>[def. 18.1]</sup> 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

$N(v)$ : The neighborhood of a vertex  $v_i \in \mathcal{V}$  is the set of all adjacent vertices:  
$$N(v_i) = \{v_k \in \mathcal{V} : \exists e_k = \{v_i, v_j\} \in \mathcal{E}, \forall v_j \in \mathcal{E}\} \quad (18.1)$$

## Degree Matrix

### Definition 18.6 Degree of a Vertex

$\delta$ : The degree of a vertex  $v$  is the cardinality of the neighborhood<sup>[def. 18.5]</sup> – 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

**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:  
$$\mathbf{D}_{i,j} := \begin{cases} \deg(v_i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (18.3)$$

## Edges

### Definition 18.8 Edges

$\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ : Represent some relation between edges<sup>[def. 18.4]</sup> 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  $\mathcal{G}$  with  $n = |\mathcal{V}|$  has between  $\left[0, \frac{1}{2}n(n-1)\right]$  edges.

## Graph Representations

### Adjacency Matrix

#### Definition 18.9 (unweighted) Adjacency Matrix

**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:  
$$\mathbf{A}_{i,j} := \begin{cases} 1 & \text{if } \exists e(i, j) \\ 0 & \text{otherwise} \end{cases} \quad (18.5)$$

#### Definition 18.10 weighted Adjacency Matrix

**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:  
$$\mathbf{A}_{i,j} := \begin{cases} \theta_{ij} & \text{if } \exists e(i, j) \\ 0 & \text{otherwise} \end{cases} \quad (18.6)$$

### Diagonal Elements

For a graph without self-loops the diagonal elements of the adjacency are all zero.

## Adjacency List

definition

pros and cons list vs matrix

## Operations on Graphs

### 1. Walks

**Definition 19.1 Walk:** A walk of a graph  $\mathcal{G}$  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  $\mathcal{G}$  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<sup>[def. 19.3]</sup> of a graph  $\mathcal{G}$  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:  
$$U \subseteq \mathcal{V} \quad \text{and} \quad F \subseteq \mathcal{E} \quad (20.1)$$

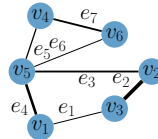
### 2. Components

**Definition 20.2 Component:** A connected component of a graph  $\mathcal{G}$  is a connected<sup>[def. 20.6]</sup> subgraph<sup>[def. 20.1]</sup> 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  $c(\mathcal{G})$ .

### 3. Weighted Graph

#### Definition 20.3 Weighted Graph:

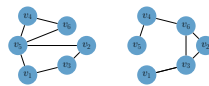
Is a graph  $\mathcal{G}$  where edges are associated with a weight:  
$$\exists \theta_i := \text{weight}(e_i) \quad \forall e_i \in \mathcal{E}$$



### 4. Spanning Graph

#### Definition 20.4 Spanning Graph:

Is a subgraph<sup>[def. 20.1]</sup>  $\mathcal{H} = (U, F)$  of a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  for which it holds:  
$$U = \mathcal{V} \quad \text{and} \quad F \subseteq \mathcal{E} \quad (20.2)$$



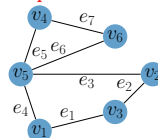
#### 4.1. Minimum Spanning Graph

**Definition 20.5 Minimum Spanning Graph:** Is a spanning graph<sup>[def. 20.4]</sup>  $\mathcal{H} = (U, F)$  of a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with minimal weights/distance of the edges.

### 5. Connected Graphs

#### Definition 20.6 (Weakly) Connected Graph:

Is a graph  $\mathcal{G}$ <sup>[def. 18.1]</sup> 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<sup>[def. 20.8]</sup> is called strongly connected if every nodes is reachable from every other node.

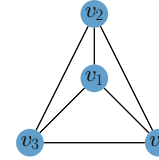
**Corollary 20.2 Components of Connected Graphs:** A connected Graph<sup>[def. 20.6]</sup> consist of one component  $c(\mathcal{G}) = 1$ .

#### 5.1. Fully Connected/Complete Graph

#### Definition 20.7 Fully Connected/Complete Graph:

Is a connected graph  $\mathcal{G}$ <sup>[def. 20.6]</sup> where each node is connected to every other node.  
$$\exists e \forall \{v_i, v_j\} \quad \forall v_i, v_j \in \mathcal{V} \quad (20.4)$$

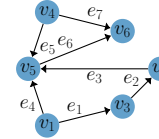
$$|\mathcal{V}| = \frac{1}{2} |\mathcal{V}| (|\mathcal{V}| - 1) \quad (20.5)$$



#### 5.2. Directed Graphs

#### Definition 20.8 Directed Graph/Digraph (DG):

A directed graph  $\mathcal{G}$  is a graph where edges are direct arcs<sup>[def. 20.9]</sup>.



**Definition 20.9 Directed Edges/Arcs:** Represent some directional relationship between edges<sup>[def. 18.4]</sup> 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

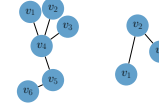
#### Definition 21.1 Acyclic Graphs:

Are graphs<sup>[def. 18.1]</sup> where no cycles<sup>[def. 19.4]</sup> exist.

## Forests

#### Definition 21.2 Forests:

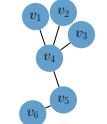
Are acyclic graphs<sup>[def. 21.1]</sup>:



## Trees

#### Definition 21.3 Trees:

Are acyclic graphs<sup>[def. 21.1]</sup> that are connected<sup>[def. 20.6]</sup>.



## 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<sup>[def. 21.5]</sup>, 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.

Red-black tree

AVL trees

#### 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<sup>[def. 21.6]</sup> where every parent is smaller/larger (min-heap/max-heap) than its children.

## Tries/Prefix Trees

#### Definition 21.10 Prefix Tree:

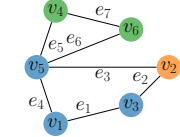
Is a tree special kind of tree where each node can have multiple children. It is usually used for prefix lookup of words, where words with the same prefix share the same nodes. It can reduce lookup time from  $\mathcal{O}(M \log N)$  for a word of size  $M$  with  $N$  total words to  $\mathcal{O}(M)$ . 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  $\mathcal{G}$  a layering of the graph is a partition of its node set  $\mathcal{V}$ <sup>[def. 18.4]</sup> into subsets

$$\begin{aligned} \{V_1, \dots, V_L\} \subseteq \mathcal{V} \\ \text{s.t. } \mathcal{V} = V_1 \cup \dots \cup V_L \end{aligned} \quad (21.2)$$



### 2. Bisection Algorithms

#### 2.1. Local Approaches

#### 2.2. Global Approaches

##### 2.2.1. Spectral Decomposition

**Definition 21.12 Graph Laplacian (Matrix)**  $\mathbf{L}(\mathcal{G})$ : Given a graph with  $n$  vertices and  $m$  edges has a graph laplacian matrix defined as:

$$\mathbf{L} = \mathbf{A} - \mathbf{D} \quad l_{ij} := \begin{cases} -1 & \text{if } i \neq j \text{ and } e_{ij} \in \mathcal{E} \\ 0 & \text{if } i \neq j \text{ and } e_{ij} \notin \mathcal{E} \\ \deg(v_i) & \text{if } i = j \end{cases} \quad (21.3)$$

**Corollary 21.2 title:**

##### 2.2.2. Inertial Bisection

## Proofs

## Model Parameter Estimation

Proof 22.1: 6.10:

$$\begin{aligned} p(\mathbf{X}, \mathbf{y}, \theta) &= \frac{p(\theta|\mathbf{X}, \mathbf{y})p(\mathbf{X}, \mathbf{y})}{p(\mathbf{y}|\mathbf{X}, \theta)p(\mathbf{X}, \theta)} \\ \frac{p(\theta|\mathbf{X}, \mathbf{y})p(\mathbf{X}, \mathbf{y})}{p(\mathbf{y}|\mathbf{X}, \theta)p(\mathbf{X}, \theta)} &= \frac{p(\theta|\mathbf{X}, \mathbf{y})p(\mathbf{y}|\mathbf{X})p(\mathbf{X})}{p(\mathbf{y}|\mathbf{X}, \theta)p(\mathbf{X}, \theta)} \\ &= p(\mathbf{y}|\mathbf{X}, \theta)p(\theta|\mathbf{X})p(\mathbf{X}) \\ &\stackrel{\text{eq. (6.6)}}{=} p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)p(\mathbf{X}) \\ \Rightarrow p(\theta|\mathbf{X}, \mathbf{y}) &= \frac{p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)p(\mathbf{X})}{p(\mathbf{y}|\mathbf{X})p(\mathbf{X})} \end{aligned}$$

### 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  $\mathbf{X}$ )

## Generative Models

1. Diffusion Models

Proof 23.1 One Step Forward Diffusion Model<sup>[def. 17.10]</sup>: Let  $\{\epsilon\}_{i=1}^T \sim \mathcal{N}(0, \mathbf{I})$ :

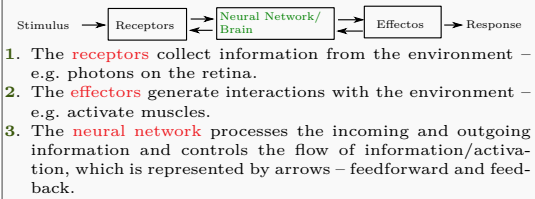
$$\begin{aligned} \mathbf{x}_t &= \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \epsilon_t = \mathbf{x}_{t-1} \sqrt{\alpha_t} + \sqrt{1 - \alpha_t} \epsilon_t \\ &= \left( \mathbf{x}_{t-2} \sqrt{\alpha_{t-1}} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-1} \right) \sqrt{\alpha_t} + \sqrt{1 - \alpha_t} \epsilon_t \\ &= \mathbf{x}_{t-2} \sqrt{\alpha_t \alpha_{t-1}} + \underbrace{\sqrt{\alpha_t (1 - \alpha_{t-1})} \epsilon_{t-1}}_{:=Y} + \underbrace{\sqrt{1 - \alpha_t} \epsilon_t}_{:=Z} \\ Y &\sim \mathcal{N}(0, \alpha_t (1 - \alpha_{t-1})) \quad Z \sim \mathcal{N}(0, 1 - \alpha_t) \\ Y + Z &\stackrel{??}{=} \mathcal{N}(0, \alpha_t (1 - \alpha_{t-1}) + (1 - \alpha_t)) \\ &\quad \mathcal{N}(0, 1 - \alpha_t \alpha_{t-1}) = \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-1} \\ \mathbf{x}_t &= \mathbf{x}_{t-2} \sqrt{\alpha_t \alpha_{t-1}} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-1} \\ &\quad \vdots \\ \mathbf{x}_t &= \mathbf{x}_{t-2} \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_0} + \sqrt{1 - \alpha_t \alpha_{t-1} \cdots \alpha_0} \epsilon_0 \\ \mathbf{x}_t &= \mathbf{x}_{t-2} \sqrt{\bar{\alpha}_t} + \sqrt{1 - \bar{\alpha}_t} \epsilon_0 \end{aligned}$$



Deep Learning Submodule

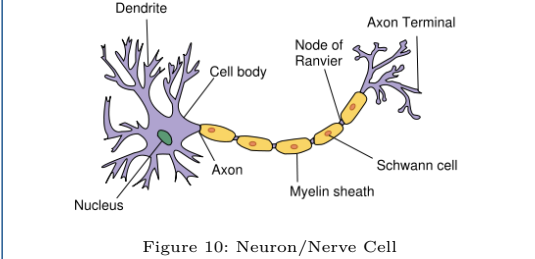
Biological Neural Networks

The human nervous system can be broken down into three stages:



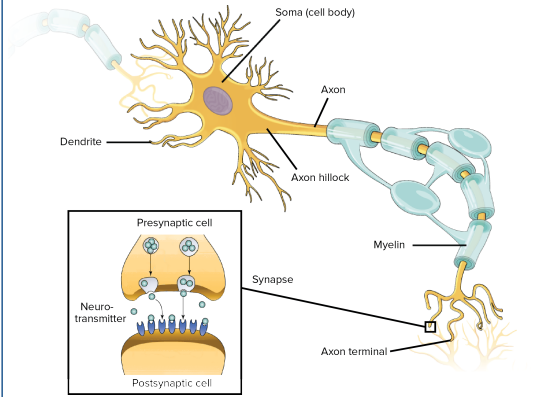
1. Neurons/Nerve Cells

**Definition 24.1 Neurons/Nerve Cells:** Are the basic building blocks of a neural networks, that processes and transmits information through electrical and chemical signals.



**Definition 24.2 Axon:** Signals are sent to other neurons via axons.

**Definition 24.3 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.



**Definition 24.4 The Soma and Dendrites:** Are responsible for receiving synaptic signals from other neurons.

**Definition 24.5 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 expiatory 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.

2.1. **Mathematical Abstraction**

We can view a neuron as real valued function that takes n-real valued inputs and maps them to true or false

$$f : \mathbb{R}^n \mapsto \{0, 1\} \tag{24.1}$$

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:

$$f : \mathbb{R}^n \times \mathbb{R}^d \mapsto \{0, 1\} \tag{24.2}$$

**Note**  
We may also view a neuron as a mapping to a real valued output, if we view the output as firing rate/frequency.

$$f : \mathbb{R}^n \times \mathbb{R}^d \mapsto \mathbb{R} \tag{24.3}$$

# History

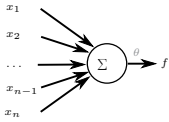
## Neurons

### 1. McCulloch-Pitts (MCP) Neuron 1943

#### Definition 25.1 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

#### 1.1. Logical Units

Any task or phenomenon that can be represented as a logic function can be modeled by a network of MP-neurons.

#### Proposition 25.1 [example 25.1]

##### Disjunctive Normal Form of MCPs:

**Given:** the tuple of the set of all activities  $\Sigma = \{\sigma_{i=1}\}^n$  and the threshold  $\theta$ .

Let  $\mathcal{I}$  be the *set* of all subsets  $I$  of  $\Sigma$  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:

$$f(\mathbf{x}; \sigma, \theta) = \bigvee_{I \in \mathcal{I}} \left( \bigwedge_{i \in I} x_i \bigwedge_{i \notin I} \neg x_i \right) \quad (25.1)$$

#### Example 25.1 NAND and DNF:

$x_1$	$x_2$	$\sum_{i=1}^2 \sigma_i x_i$	$f$	In order to satisfy NAND we need to set $\sigma_1 = -1$ , $\sigma_2 = -1$ and $\theta = -2$ . $\mathcal{I} = \{\emptyset, \{1\}, \{2\}\}$
0	0	0	1	
0	1	-1	1	
1	0	-1	1	
1	1	-2	0	

$f(\mathbf{x}; \sigma, \theta) = (\neg x_1 \wedge \neg x_2) \vee (x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)$

### 2. Turing Machines 1948

**Definition 25.2 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 25.3 Type A Machines:

Are networks of NAND units, operating in a clocked manner:

$$y(t+1) = 1 - x_1(t)x_2(t)$$
$$x_1(t), x_2(t) \in \{0, 1\} \quad \forall t \quad (25.2)$$

$x_1$	$x_2$	$y$
0	0	1
1	1	1
0	0	1
1	1	0

#### Note

Alan turing used NAND as all boolean function's can be decomposed into NAND-units.

Finish/understand Type B machines

#### Definition 25.4 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 25.5 r-Sparse Boolean Vectors:

### 4. The Perceptron

1958+

**Definition 25.6 Perceptron Classifier/Unit:** The perceptron is a *linear classifier*<sup>[def. 4.20]</sup> that uses the *threshold function* in order to classify dichotomies<sup>[def. 4.21]</sup>:

$$\hat{y}(\mathbf{x}, \mathbf{w}) = \text{sign}(\mathbf{x}^\top \mathbf{w}) = \begin{cases} +1 & \text{if } \sum_{i=1}^d w_i x_i \geq 0 \\ -1 & \text{else} \end{cases} \quad (25.3)$$

#### Notes

- The *activation function* of the perceptron corresponds to the heaviside step function??
- The perceptron unit is also termed *single-layer-perceptron*, to distinguish it from the misnomer for the multilayer perceptron.

#### 4.1. Number of Errors

##### 4.1.1. Existence of a Separating Solution

separable solution

##### 4.1.2. Rate of Convergence

#### Definition 25.7 Novikovs Theorem:

##### 4.1.3. Uniqueness

covers theorem probably best to linear classifier

### 5. Learning

#### 5.1. Hebb's Rule

1949

# Linear Networks

## 1. Autoencoder

### Definition 26.1 Autoencoder

Let  $X = \{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^d$  and let  $\mathbf{X}^\top = [\mathbf{x}_1 \dots \dots \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  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  $\hat{\mathbf{x}} = (\psi \circ \phi)\mathbf{x}$  s.t.

$$\begin{aligned} \theta = \{\phi, \psi\} &= \arg \min_{\phi, \psi} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - (\psi \circ \phi)\mathbf{x}_i\|^2 \\ &= \arg \min_{\phi, \psi} \frac{1}{2n} \|\mathbf{X}^\top - (\psi \circ \phi)\mathbf{X}^\top\|_F^2 \end{aligned} \quad (26.1)$$

**Explanation 26.1** (Definition 26.1). An autoencoder is a neural network that is supposed to learn the identity map by first encoding and then decoding the input.

### 1.1. Linear Autoencoder

**Definition 26.2 Linear Autoencoder**  $\hat{\mathbf{X}} = \mathbf{DCX}$ : Let  $X = \{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^d$  and let  $\mathbf{X}^\top = [\mathbf{x}_1 \dots \dots \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  be the *transposed design matrix*. An auto encoder is defined as:

**Encoder**  $\mathbf{C} \in \mathbb{R}^{p \times d}: \mathbb{R}^d \mapsto \mathbb{R}^p$  s.t.  $\hat{\mathbf{x}} = \mathbf{DCx}$   
**Decoder**  $\mathbf{D} \in \mathbb{R}^{p \times d}: \mathbb{R}^p \mapsto \mathbb{R}^d$

$$\begin{aligned} \theta = \{\mathbf{C}, \mathbf{D}\} &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{DCx}_i\|^2 \\ &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \|\mathbf{X}^\top - \mathbf{DCX}^\top\|_F^2 \end{aligned} \quad (26.3)$$

#### Note

Equation (27.2) is a degenerated problem and has single no-global minima but many minima that are optimal.

**Corollary 26.1 Non-unique solution** [proof 37.15]: The solution  $\mathbf{C}^*$ ,  $\mathbf{D}^*$  are *global* degenerative/non-unique solutions.

#### 1.1.1. Loss Function

add/find derivative of lecture slides

exercise 10/11 2019

get familiar with complex matrix derivatives

**Proposition 26.1 Rank of  $\hat{\mathbf{X}}$ :** As a result of the matrix multiplication and ?? the rank of  $\hat{\mathbf{X}}$  cannot increase:  
 $\text{rank}(\hat{\mathbf{X}}) \leq p$  (26.4)

**Proposition 26.2 Optimal Solution** [proof 37.16]: Due to ?? the best solution the linear autencoder can find is given by:

$$\mathbf{C}^* = \mathbf{U}_k^\top \quad \text{and} \quad \mathbf{D}^* = \mathbf{U}_k \quad (26.5)$$

maybe add comparison/similarity to PCA

maybe add deep linear stuff

<https://courses.media.mit.edu/2019fall/naa022/whiten.pdf>

## 2. Variational Autoencoder

**Definition 26.3 Linear Autoencoder**  $\hat{\mathbf{X}} = \mathbf{DCX}$ : Let  $X = \{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^d$  and let  $\mathbf{X}^\top = [\mathbf{x}_1 \dots \dots \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  be the *transposed design matrix*. An auto encoder is defined as:

**Encoder**  $\mathbf{C} \in \mathbb{R}^{p \times d}: \mathbb{R}^d \mapsto \mathbb{R}^p$  s.t.  $\hat{\mathbf{x}} = \mathbf{DCx}$   
**Decoder**  $\mathbf{D} \in \mathbb{R}^{p \times d}: \mathbb{R}^p \mapsto \mathbb{R}^d$

$$\begin{aligned} \theta = \{\mathbf{C}, \mathbf{D}\} &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{DCx}_i\|^2 \\ &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \|\mathbf{X}^\top - \mathbf{DCX}^\top\|_F^2 \end{aligned} \quad (26.6)$$

#### Note

Equation (27.2) is a degenerated problem and has single no-global minima but many minima that are optimal.

**Corollary 26.2 Non-unique solution** [proof 37.15]: The solution  $\mathbf{C}^*$ ,  $\mathbf{D}^*$  are *global* degenerative/non-unique solutions.

#### 2.0.1. Loss Function

add/find derivative of lecture slides

exercise 10/11 2019

get familiar with complex matrix derivatives

**Proposition 26.3 Rank of  $\hat{\mathbf{X}}$ :** As a result of the matrix multiplication and ?? the rank of  $\hat{\mathbf{X}}$  cannot increase:  
 $\text{rank}(\hat{\mathbf{X}}) \leq p$  (26.8)

**Proposition 26.4 Optimal Solution** [proof 37.16]: Due to ?? the best solution the linear autencoder can find is given by:

$$\mathbf{C}^* = \mathbf{U}_k^\top \quad \text{and} \quad \mathbf{D}^* = \mathbf{U}_k \quad (26.9)$$

## Generative Models

### Variational Autoencoder

**Definition 27.1 Linear Autoencoder**  $\hat{\mathbf{X}} = \mathbf{DCX}$ : Let  $X = \{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^d$  and let  $\mathbf{X}^\top = [\mathbf{x}_1 \dots \dots \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  be the *transposed design matrix*. An auto encoder is defined as:

**Encoder**  $\mathbf{C} \in \mathbb{R}^{p \times d}: \mathbb{R}^d \mapsto \mathbb{R}^p$  s.t.  $\hat{\mathbf{x}} = \mathbf{DCx}$   
**Decoder**  $\mathbf{D} \in \mathbb{R}^{p \times d}: \mathbb{R}^p \mapsto \mathbb{R}^d$

$$\begin{aligned} \theta = \{\mathbf{C}, \mathbf{D}\} &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{DCx}_i\|^2 \\ &= \arg \min_{\mathbf{C}, \mathbf{D}} \frac{1}{2n} \|\mathbf{X}^\top - \mathbf{DCX}^\top\|_F^2 \end{aligned} \quad (27.2)$$

#### Note

Equation (27.2) is a degenerated problem and has single no-global minima but many minima that are optimal.

**Corollary 27.1 Non-unique solution** [proof 37.15]: The solution  $\mathbf{C}^*$ ,  $\mathbf{D}^*$  are *global* degenerative/non-unique solutions.

#### 0.0.1. Loss Function

add/find derivative of lecture slides

exercise 10/11 2019

get familiar with complex matrix derivatives

**Proposition 27.1 Rank of  $\hat{\mathbf{X}}$ :** As a result of the matrix multiplication and ?? the rank of  $\hat{\mathbf{X}}$  cannot increase:  
 $\text{rank}(\hat{\mathbf{X}}) \leq p$  (27.3)

**Proposition 27.2 Optimal Solution** [proof 37.16]: Due to ?? the best solution the linear autencoder can find is given by:

$$\mathbf{C}^* = \mathbf{U}_k^\top \quad \text{and} \quad \mathbf{D}^* = \mathbf{U}_k \quad (27.4)$$

## Generative Adversarial Networks (GANs)

### Definition 28.1 Gan:

$$\min_G \max_D \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{x \sim p_Z(z)} [1 - D(\log D(x))] \quad (28.1)$$

### 1. Proofs

Proof 28.1 [Corollary 27.1]: Let  $\mathbf{C}^*$ ,  $\mathbf{D}^*$  be the global solution of problem Equation (27.2). Then we may define a matrix  $\mathbf{A}$  s.t.:

$$\mathbf{D}^* \mathbf{C}^* = \mathbf{D}^* (\mathbf{A}^{-1} \mathbf{A}) \mathbf{C}^* = \tilde{\mathbf{D}}^* \tilde{\mathbf{C}}^* \quad (28.2)$$

Proof 28.2 Proposition 27.2:

$$\begin{aligned} \mathbf{DCX} &\stackrel{??}{=} \mathbf{DCU} \Sigma \mathbf{V}^H = \mathbf{U}_k^\top \underbrace{\mathbf{U}_k \mathbf{U}}_{\stackrel{??}{=} \mathbf{I}_m} \Sigma \mathbf{V}^H \\ &= \mathbf{U}_k \Sigma_m \mathbf{V}_m^H = \mathbf{X}_m \end{aligned}$$

## Units/Activation Functions

**Definition 29.1 Activation Pattern:** Is the set of neurons of a neural network that are active/inactive for a given input  $\mathbf{x}$ .

**Definition 29.2 Saturation:** Corresponds to input value  $z$  that causes the maximal activity of the unit (positive or negative).

**Proposition 29.1 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.

### 1. Threshold Units

#### 1.1. Linear Units

Has no maximum saturation rate and may hence grow to infinity or minus infinity.

$$\varphi: \mathbb{R}^n \mapsto \mathbb{R} \quad \varphi(\mathbf{x}) = \sum_{i=1}^N w_i x_i + b = \mathbf{w}^T \mathbf{x} + b$$

$\varphi(\mathbf{x}) = \mathbf{x}$  Activation function is the identity

#### 1.2. 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).

$$\varphi: \mathbb{R}^n \mapsto \{0, 1\} \quad \varphi(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^N w_i x_i + b \right) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

**Problem:** p.w. continuous constant  $\Rightarrow$  zero gradient.

### 2. Sigmoid Units

**Definition 29.3 Sigmoid Functions:** Are functions whose graph is a characteristic “S”-shaped curve.

#### 2.1. Sigmoid/Logistic Units

**Definition 29.4 Sigmoid/Logistic Unit**  $\sigma: \mathbb{R} \mapsto [0, 1]$ : Is a smooth non-linear function that can not have negative activities:

$$\sigma(\mathbf{x}^T \theta) = \frac{1}{1 + e^{-\mathbf{x}^T \theta}} \quad \sigma'(z) = \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z) \quad \lim_{z \rightarrow \pm\infty} \sigma'(z) = 0$$

(29.1)

**Explanation 29.1** (Sigmoid/Logistic Function).

$$\sigma(\mathbf{x}^T \theta) = \begin{cases} 0 & -\mathbf{x}^T \theta \text{ large} \\ 1 & \text{if } \mathbf{x}^T \theta \text{ large} \\ 0.5 & \mathbf{x}^T \theta = 0 \end{cases}$$

**Corollary 29.1 Smooth Function:**

$\sigma$  is a smooth function  $\sigma \in C^\infty$ .

(proof recursion formula)

**Corollary 29.2 Rotational Symmetry** around  $(0, 1/2)$ :  $1 - \sigma(z) = \sigma(-z) \iff \sigma(-z) = 1 - \sigma(z)$  (29.2)

**Property 29.1 Inverse**

$$\sigma^{-1}(z) = \ln \left( \frac{t}{1-t} \right) \quad (29.3)$$

#### Pros

- Great for classification problems.

#### Cons

- Saturated logistic units are non-informative see Corollary 29.3.

#### Note

- Useful for the output layer in deep neural networks.
- Proven to work great for three layer and two layer neural networks, particularly classification problems. Notice the hill-shaped derivative of the function which pushes the network to “move down the hill” to either side, giving the network more distinction when classifying.

#### 2.1.1. Saturated Logistic Units

**Corollary 29.3 Saturated logistic units:**

A saturated logistic unit<sup>[def. 29.4]</sup> (large input values  $z$ ) has a zero/vanishing gradient<sup>[def. 30.11]</sup>:

$$\lim_{z \rightarrow \pm\infty} \sigma'(z) = \lim_{z \rightarrow \pm\infty} \sigma(z)\sigma(-z) = 0 \quad (29.4)$$

$\Rightarrow$  parameters will not be updated as the gradient/ $\epsilon$  is zero.

add switch activation

#### 2.2. Hyperbolic Tangent

**Definition 29.5**

**Hyperbolic Tangent/Tanh Unit**  $\tanh: \mathbb{R} \mapsto [-1, 1]$ : Is a smooth symmetric non-linear function:

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{\sinh(z)}{\cosh(z)} \quad \tanh'(z) = 1 - \tanh^2(z)$$

(29.5) (29.6)

#### Property

**Property 29.2**

**Connection to Logistic Unit:**  $\tanh(z) = 2\sigma(2z) - 1$  [proof 37.3]

#### 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.

#### Notes

Works almost always better than the sigmoid  $\sigma$  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.

## Softmax – Generalized Sigmoid Unit

**Definition 29.6**

**Softmax Function**

Is a non-negative, non-linear, normalized function that generalizes the logistic unit<sup>[def. 29.4]</sup> to multiple classes:  $\sigma^{\max}: \mathbb{R}^c \mapsto \mathbb{R}^c$ :  $\sigma_k^{\max}(\mathbf{x}; \Theta) = [\sigma_1^{\max}, \dots, \sigma_c^{\max}]$   $\sigma_k^{\max} \in [0, 1] \quad \forall k$

$$\sigma_i^{\max}(\mathbf{x}; \Theta) = \frac{\exp(\mathbf{x}^T \theta_i)}{\sum_{j=1}^c \exp(\mathbf{x}^T \theta_j)} \quad \Theta = [\theta_1, \dots, \theta_c] \quad (29.7)$$

$$\|\sigma^{\max}\|_1 = \sum_{j=1}^c \sigma_j = 1 \quad \sigma_i > 0 \quad \forall i = 1, \dots, c$$

$$\nabla_{z_j} \sigma_i^{\max}(z) = \begin{cases} \sigma_i^{\max}(z)(1 - \sigma_i^{\max}(z)) & i = j \\ -\sigma_i^{\max}(z)\sigma_j^{\max}(z) & i \neq j \end{cases} = \sigma_i^{\max}(z)(\delta_{ij} - \sigma_j^{\max}(z)) \quad (29.8)$$

Logits  $\mathbf{x}$

$$\begin{bmatrix} 3.1 \\ -12.2 \\ \vdots \\ 3 \end{bmatrix}$$

Softmax

$$\frac{e^{\mathbf{x}^T \theta_i}}{\sum_{j=1}^c e^{\mathbf{x}^T \theta_j}}$$

Probabilities

$$\begin{bmatrix} 0.42 \\ 0.04 \\ \vdots \\ 0.3 \end{bmatrix}$$

**Explanation 29.2.**

- The exponential function makes sure that we cannot get negative probabilities.
- If there exist logits<sup>[def. 29.7]</sup> 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 they other pushed towards zero.

**Corollary 29.4 Relationship to Sigmoid Unit**[proof 37.6]:

$$\text{soft-max with } \theta_1, \theta_2 \iff \text{sigmoid with } \theta := \theta_1 - \theta_2 \quad (29.9)$$

**Corollary 29.5 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:  $\theta_i \leftarrow \theta_i - \theta_k \quad i = 1, \dots, c \quad \text{s.t.} \quad \theta_k \leftarrow 0 \quad (29.10)$

#### 2.2.1. Logits

**Definition 29.7 Logit(s)**

Is the logarithm of the odds of an event and equals the inverse of the sigmoid function:

$$\text{logit}(p) = \ln \frac{p}{1-p} = \ln(\text{odds}(p)) \quad (29.11)$$

**Corollary 29.6 Logits and Neural Networks:** The un-normalized and possibly negative inputs to the softmax function are thus called logits.

#### 2.2.2. 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 29.8**

**Numerical Stable Softmax:** Softmax function that avoid over and underflow:  $\sigma^{\max}(\mathbf{z}) = [\sigma_1^{\max}, \dots, \sigma_c^{\max}]$   $\sigma_k^{\max} \in [0, 1] \quad \forall k$

$$\sigma_i^{\max}(\mathbf{z}) = \frac{\exp(\mathbf{z}_i + a)}{\sum_{j=1}^c \exp(\mathbf{z}_j + a)} \quad a = -\max(z_1, \dots, z_c) \quad (29.12)$$

**Explanation 29.3.** • 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.

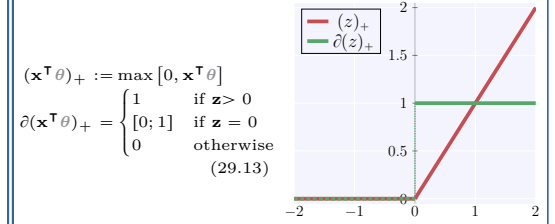
### 3. Rectified Units

**Definition 29.9 Rectified Units:** Is the use of piece wise linear units  $\sigma_{pw}^0$ .

#### 3.1. Rectified Linear Units

**Definition 29.10**

**Rectified Linear Units (ReLU)**



**Explanation 29.4.** ReLU splits the input space into two half spaces:

$$\mathcal{H}_\theta^+ = \{\mathbf{x} : \theta^T \mathbf{x} > 0\} \quad \text{and} \quad \mathcal{H}_\theta^- = \{\mathbf{x} : \theta^T \mathbf{x} < 0\}$$

separated by the hyperplane  $\mathcal{H}_\theta^0$  and constant on  $\mathcal{H}_\theta^0 \cup \mathcal{H}_\theta^-$

**Corollary 29.7 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:

$$\text{If } z_{lj} = 0 \implies \partial(z_{lj})_+ = 0 \quad (29.14)$$

$$\implies \nabla_{\theta_{lj}} = \nabla_{z_{k-1}} z_{lj} = 0 \quad (29.15)$$

Thus if  $z_{lj}(\mathbf{x}^t) = 0$  for all inputs of a training sample

- $\mathbf{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^T\}$  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 29.8 Activation Pattern:** The activation pattern<sup>[def. 29.1]</sup> of a ReLU layer can be described by a heavy-side function??:

$$H(\Theta \mathbf{x}) \in \{0, 1\}^m \quad \mathbf{x} \in \mathbb{R}^n \quad (29.16)$$

as each unit will either be active or inactive.

**Corollary 29.9 Jaccobi Matrix:** The Jaccobi matrix?? of a ReLU layer is sparse:

$$\partial \mathbf{F}_l := \begin{bmatrix} \partial_{\theta_{l1}}^T \\ \partial_{\theta_{l2}}^T \\ \vdots \\ \partial_{\theta_{l1}}^T \end{bmatrix}, \quad \tilde{\theta}_{lj} = \begin{cases} 0 & \text{if } z_{lj} = 0 \\ \partial_{lj} & \text{else} \end{cases} \quad (29.17)$$

Thus if all units are active we recover a linear network otherwise inactive units are “pruned away”.

#### Pros

- Non-linear gradient  $\Rightarrow$  easy to optimize.

#### Cons

- No gradient information for negative entries.
- Dying ReLU proble

### 3.2. Absolute Value Units

**Definition 29.11** [proof ??]  
**Absolute Value (AbsU)**

$|z| := \begin{cases} z & \text{if } z \geq 0 \\ -z & \text{otherwise} \end{cases}$

$\partial(|z|) = \begin{cases} 1 & \text{if } z > 0 \\ [0, 1] & \text{if } z = 0 \\ 0 & \text{otherwise} \end{cases}$

(29.18)

**Corollary 29.10** [proof 37.9]  
**Connection to Absolute value Unit:**

$(z)_+ = \frac{z + |z|}{2} \iff |z| = 2(z)_+ - z$  (29.19)

### 3.3. Smooth ReLU Approximations

Are no-longer “real” REU units<sup>[def. 29.9]</sup> but *smooth* approximations of the the *p.w. RELU*<sup>[def. 29.10]</sup> unit and try to combine the advantages of *rectification* and *smoothness*.

#### 3.3.1. Softplus

**Definition 29.12** **Softplus**

$\zeta : \mathbb{R} \mapsto (0; \infty):$

$\zeta(\mathbf{x}; \theta) = \ln(1 + \exp(\mathbf{x}^\top \theta))$

$\zeta'(\mathbf{x}^\top \theta) = \sigma(\mathbf{x}^\top \theta) = \frac{1}{1 + e^{-\mathbf{x}^\top \theta}}$

(29.20)

#### 3.3.2. Exponential Linear Unit

**Definition 29.13**  
**Exponential Linear Unit ELU**

$\varphi(\mathbf{z}) = \begin{cases} \mathbf{z} & \text{if } \mathbf{z} \geq 0 \\ \exp(\mathbf{z}) - 1 & \text{else} \end{cases}$

(29.21)

$\frac{d}{d\mathbf{z}} \varphi(\mathbf{z}) = \begin{cases} 1 & \text{if } \mathbf{z} > 0 \\ \exp(\mathbf{z}) & \text{otherwise} \end{cases}$

(29.22)

### 3.4. Leaky ReLU

**Definition 29.14** **Leaky ReLU**

$\varphi(\mathbf{z}) := \begin{cases} \mathbf{z} & \text{if } \mathbf{z} \geq 0 \\ \epsilon \mathbf{z} & \text{otherwise} \end{cases}$

(29.23)

$\partial(\varphi(\mathbf{z})) = \begin{cases} 1 & \text{if } \mathbf{z} > 0 \\ [\epsilon, 1] & \text{if } \mathbf{z} = 0 \\ \epsilon & \text{else} \end{cases}$

(29.24)

### 3.5. Approximation Power of pw. Lin. Functions

#### 3.5.1. Hinge Functions

**Definition 29.15** **Hinge Functions:**

Are two hyperplanes?? that are cut in half and glued together at their intersection:

$g(x) = \max(\theta_1^\top \mathbf{x} + b_1, \theta_2^\top \mathbf{x} + b_2)$

(29.25)

**Corollary 29.11** **Intersection:** The intersection is given by:  
 $H_1 = H_2 \iff H_2 - H_1 = 0 \iff \theta_2^\top \mathbf{x} + b_2 - \theta_1^\top \mathbf{x} + b_1 = 0$   
 $\Rightarrow (\theta_2 - \theta_1)^\top \mathbf{x} + (b_2 - b_1) =: \tilde{\theta} + \tilde{b} = 0$

**Corollary 29.12** **Representation:** Hinge functions<sup>[def. 29.15]</sup> can be represented as combination of a linear part?? and an absolute value or relu unit:

$\max(f, g) = \frac{1}{2}(f + g + |f - g|)$  (29.26)

#### Note

Absolute value functions in turn can be represented by relu functions.

#### 3.5.2. k-Hinge Functions–Maxout Units Polyhedra

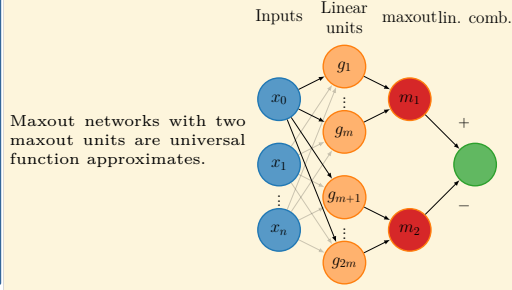
**Theorem 29.1** **Hinge Representation:**  
Polyhedra  $f \in \mathbb{R}^n$  can be represented by (finite) *k*-Hinge functions:  
 $\exists \mathcal{A} \subset \mathbb{R}^{n+1} : (\theta, b) \in \mathcal{A} \text{ and } f(x) = \max_{(\theta, b) \in \mathcal{A}} \{\theta^\top \mathbf{x} + b\}$

(29.27)

**Theorem 29.2** (Wang, 2004):  
Every continuous piece wise linear function  $f \in C_{pw}^0$  can be written as the difference of *two*

#### Maxout Units/*k*-Hinge Functions

**Theorem 29.3** **Minimal Non-Linearity:**



## 4. Proofs

**Definition 29.4**

$\sigma'(z) = \frac{\partial}{\partial z} (1 + e^{-z})^{-1} = -(1 + e^{-z})^{-2} \frac{\partial}{\partial z} (1 + e^{-z})$

Proof 29.1:

$= -\frac{1}{(1 + e^{-z})^2} e^{-z} (-1) = \frac{1}{(1 + e^{-z})^2} (1 - e^{-z})$

$= \frac{1}{(1 + e^{-z})} - \frac{1}{(1 + e^{-z})^2} = \sigma(z) - \sigma(z)^2$

$\sigma'(z) = \frac{\partial}{\partial z} (1 + e^{-z})^{-1} = -(1 + e^{-z})^{-2} e^{-z} (-1)$

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

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

Proof 29.2: Property 29.1

$$\sigma\left(\ln\left(\frac{t}{1-t}\right)\right) = \frac{1}{\frac{1}{1-t}} = t \quad (29.28)$$

Proof 29.3: Property 29.2

$$2\sigma(2z) - 1 = 2 \frac{1}{1 + e^{-2z}} - 1 = \frac{2}{1 + e^{-2z}} - \frac{1 + e^{-2z}}{1 + e^{-2z}} = \frac{1 - e^{-2z}}{1 + e^{-2z}} = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \tanh(z)$$

Proof 29.4: Definition 29.5

$$\tanh(z) = \frac{\frac{\partial}{\partial z} \sinh(z)}{\frac{\partial}{\partial z} \cosh(z)} = \frac{\frac{\partial}{\partial z} \sinh(z) \cdot \cosh(z) - \frac{\partial}{\partial z} \cosh(z) \cdot \sinh(z)}{\cosh^2(z)}$$

$$= \frac{\cosh^2(z) - \sinh^2(z)}{\cosh^2(z)} = 1 - \frac{\sinh^2(z)}{\cosh^2(z)}$$

Proof 29.5: Definition 29.6  
Softmax function derivative  $\nabla_{z_j} \sigma_i^{\max}(z)$ :

$i \neq j : \nabla_{z_j} \frac{\exp(\mathbf{x}^\top \theta_i)}{\sum_{k=1}^c \exp(\mathbf{x}^\top \theta_k)} =$

Q.R.  $0 - e^{x_j} e^{x_i} = -\frac{e^{x_j}}{\sum_{k=1}^c e^{x_k}} \cdot \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}}$

$= -\sigma_i^{\max}(z) \sigma_j^{\max}(z)$

$i = j : \nabla_{z_i} \frac{\exp(\mathbf{x}^\top \theta_i)}{\sum_{k=1}^c \exp(\mathbf{x}^\top \theta_k)} =$

$= \frac{e^{x_i} \sum_{k=1}^c e^{x_k} - e^{x_i} e^{x_i}}{\left(\sum_{k=1}^c e^{x_k}\right)^2}$

$= \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} - \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} \cdot \frac{e^{x_i}}{e^{x_k}}$

$= \sigma_i^{\max}(z) - (\sigma_i^{\max}(z))^2$

Proof 29.6: Corollary 37.6

$$\sigma_1^{\max}(\mathbf{x}) = \frac{e^{\mathbf{x}^\top \theta_1}}{e^{\mathbf{x}^\top \theta_1} + e^{\mathbf{x}^\top \theta_2}} = \frac{e^{-\mathbf{x}^\top \theta_1}}{e^{-\mathbf{x}^\top \theta_1} + e^{\mathbf{x}^\top \theta_2 - \mathbf{x}^\top \theta_1}} = \frac{1}{1 + e^{\mathbf{x}^\top (\theta_2 - \theta_1)}}$$

Proof 29.7: Definition 29.7

$$y = \frac{e^x}{e^x + 1 + e^{-x}} = \frac{e^x}{1 + e^x} = \frac{1 + e^x - 1}{1 + e^x} = 1 - \frac{1}{1 + e^x}$$

$\Rightarrow \frac{1}{1 + e^x} = 1 - y$

$$1 + e^x = \frac{1}{1 - y} = \frac{1 - y + y}{1 - y} = 1 + \frac{y}{1 - y}$$

$$e^x = \frac{y}{1 - y} \implies x = \ln \frac{y}{1 - y}$$

Proof 29.8: Corollary 29.10

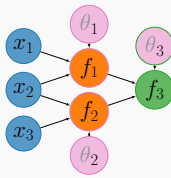
If  $z \geq 0 \implies (z)_+ = 2z - z = z$  if  $z < 0$  (29.29)

$\implies (z)_+ = 0$  (29.30)

# Neural Network Representation

## 1. DAG Representation

**Idea:** neural networks that can be represented by acyclic<sup>[def. 21.1]</sup> directed<sup>[def. 20.8]</sup> graphs<sup>[def. 18.1]</sup>.



$$\mathbf{x} \mapsto f_3(f_1(x_1, x_2, \theta_1), f_2(x_2, x_3, \theta_2), \theta_3)$$

add different version i.e. torch vs tensorflow

## 2. Layered Neural Networks

### 2.0.1. Neural Network Layers

**Definition 30.1 Layered Neural Networks:**

Are neural networks where the graph<sup>[def. 18.1]</sup> is partitioned into  $L$  layers<sup>[def. 30.2]</sup>.

**Definition 30.2 Neural Network Layers**  $1, \dots, L$ :

Are grouped partitions of nodes<sup>[def. 21.11]</sup> of the neural network.

### 2.1. Composition, Activations and Transfer Maps

**Definition 30.3 Transfer Map:**

Is a function  $\mathbf{F}^l$  that maps the input of one layer  $l-1$  onto the next layer  $l$ :

$$\mathbf{F}^l : \mathbb{R}^{n_l} \times n_{l-1} \mapsto \mathbb{R}^{n_l} \iff \mathbf{F}^l := \varphi \circ \mathbf{F}^l \quad (30.1)$$

$$\text{with} \quad \bar{\mathbf{F}}^l(\mathbf{x}) = \theta^l \mathbf{x} + \mathbf{b}^l \quad \mathbf{b}^l \in \mathbb{R}^{m^l}$$

**Definition 30.4**

**Composition of Layered Neural Networks:**

Is the composition of the  $L$  transfer maps<sup>[def. 30.3]</sup>  $\{\mathbf{F}_l\}_{l=1}^L$  w.r.t. the input, whereas the parameters get concatenated:

$$\mathbf{F} = \mathbf{F}_{L:1} = \mathbf{F}_L \circ \dots \circ \mathbf{F}_1 \quad \mathbf{F}_l : \mathbb{R}^{n_l} \mapsto \mathbb{R}^{n_{l+1}} \quad (30.2)$$

with  $\mathbf{F}(\mathbf{x}; \ominus) = \mathbf{F}_L(\mathbf{F}_{L-1}(\dots \mathbf{F}_1(\mathbf{x}; \theta_1, \dots, \theta_{L-1})\theta_L))$

$$\mathbf{z}_l := \mathbf{F}_{l:1} = (\mathbf{F}_l \circ \mathbf{F}_{l-1:1})(\mathbf{x}) = \mathbf{F}_l(\mathbf{z}_{l-1})$$

$$\mathbf{z}_0 = \mathbf{x} \quad \mathbf{z}_L = \hat{\mathbf{y}}$$

### 2.1.1. Width of a Layer

**Definition 30.5 With of a Layer:**

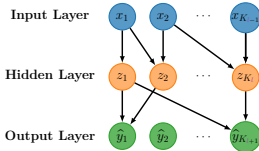
$$\text{width}_l = n_l := \dim(\text{codom}(\mathbf{F}_l)) \quad (30.3)$$

## 3. Feed Forward Neural Networks

**Definition 30.6**

**Feed Forward Neural Networks (FFNN):**

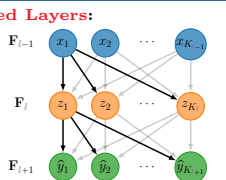
Are neural networks that can be represented by a acyclic<sup>[def. 21.1]</sup> directed<sup>[def. 20.8]</sup> graphs<sup>[def. 18.1]</sup> where information only flows into one direction i.e. along Paths<sup>[def. 19.3]</sup>.



## 4. Fully Connected Neural Network

**Definition 30.7 Fully Connected Layers:**

Are neural network layers<sup>[def. 30.2]</sup>  $\mathbf{F}_l$  whose edges are fully connected<sup>[def. 20.7]</sup> to all edges of the previous  $\mathbf{F}_{l-1}$  and the next layer  $\mathbf{F}_{l+1}$ .



### 4.1. Fully Connected FFNN

**Definition 30.8**

**(FC-FFNN)**

**Fully Connected Feed Forward Neural Networks :**

Are feed forward neural<sup>[def. 30.6]</sup> networks that consists only off fully connected layers<sup>[def. 30.7]</sup>.

#### 4.1.1. Units

**Definition 30.9 Units of a Neural Network:**

Is the sum of the widths<sup>[def. 30.5]</sup> of all layers:

$$\# \text{units} = \sum_{l=1}^L \text{width}_l \quad (30.4)$$

### 4.2. Forward Propagation

### 4.3. Backpropagation

#### 4.3.1. Derivative of the Network

**Definition 30.10**

[proof 30.1]

**Chain Rule for Neural Networks:**

$$\begin{aligned} \partial \mathbf{F} &= \prod_{l=L}^1 \partial \mathbf{F}_l \circ \mathbf{F}_{l-1:1} = \prod_{l=L}^1 \mathbf{J}_{\mathbf{F}_l} \circ \mathbf{F}_{l-1:1} \\ &\equiv \partial \mathbf{F}(\mathbf{x}) = \prod_{l=L}^1 \partial \mathbf{F}_l(\mathbf{z}_{l-1}) \end{aligned} \quad (30.5)$$

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 30.11 Vanishing Gradient Problem:**

fix/update section

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 30.1:

$$\begin{aligned} \partial \mathbf{F} &= (\partial \mathbf{F}_l \circ \mathbf{F}_{l-1:1}) \cdot \partial \mathbf{F}_{l-1:1} \\ &= (\partial \mathbf{F}_l \circ \mathbf{F}_{l-1:1}) \cdot (\partial \mathbf{F}_{l-1} \circ \mathbf{F}_{l-2:1}) \cdot \partial \mathbf{F}_{l-2:1} = \dots \end{aligned}$$

Optimization

Loss Functions

- 5.1. Zero-One/Classification Loss
- 5.2. Hinge/Perceptron Loss
- 5.3. Probabilistic-Log Likelihood Loss

**Definition 30.12 (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 does not really matter as the only vary by a constant factor  $\log_2 n = \frac{\ln 2}{\ln n}$ .

**Categorical-Log Loss**  
**Definition 30.13** [proof 37.10]  
**Log/Categorical Loss:**  
Given one-hot<sup>[def. 4.26]</sup> encoded hard labels<sup>[def. 4.27]</sup> with  $K$  individual classes  $\{c_1, \dots, c_K\}$ .  
A categorical output distribution?? induces a negative log-likelihood loss<sup>[def. 6.4]</sup> of:

$$l_y(\hat{\mathbf{y}}) = -\ln f(y|\hat{\mathbf{y}}; \Theta) = -\delta_{[y=c_k]} \cdot \ln \hat{p}_k = -\sum_{k=1}^K y_k \cdot \ln \hat{p}_k$$

(30.6)

**Explanation 30.1.** [example 38.1]  
Given a true observation/label  $y = c_j$  the categorical loss penalizes only the predictive probability  $\hat{y}_j \in \hat{\mathbf{p}} = [\hat{y}_1, \dots, \hat{y}_K]^T$  of the correct  $c_j$  by the amount it did not match 1.0.

**Cross Entropy Loss**  
**Definition 30.14** [proof 37.12]  
**Cross Entropy Loss:**  
Given soft or hard probabilistic targets<sup>[def. 4.28]</sup>  $\mathbf{p} \in \mathcal{Y}$  with  $K$  individual classes  $\mathcal{Y} = \{c_1, \dots, c_K\}$  and probabilistic predictions  $\hat{\mathbf{p}}$  the cross entropy<sup>[def. 3.7]</sup> loss is given by:

$$\text{CE} = l_{\mathbf{p}}(\hat{\mathbf{p}}) = H(\mathbf{p}, \hat{\mathbf{p}}) = -\mathbb{E}_{\mathbf{p}}[\ln \hat{\mathbf{p}}] = H(\mathbf{p}) + \text{KL}(\mathbf{p} \parallel \hat{\mathbf{p}})$$
$$\stackrel{\text{disc.}}{=} -\sum_{k=1}^K p_k \ln \hat{p}_k \tag{30.7}$$
$$\frac{\partial}{\partial z_k} l_{\mathbf{p}}(\hat{\mathbf{p}}(\mathbf{z})) = -\sum_{k=1}^K p_k \frac{1}{\hat{p}(z_k)} \frac{\partial}{\partial z_k} \hat{p}(z_k) \tag{30.8}$$

**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:

$$\hat{\mathbf{y}} = \arg \max_{k \in \{1, \dots, K\}} \hat{y}_k \tag{30.9}$$

**Attention: Log-loss vs. Cross Entropy Loss:**  
The *cross-entropy loss*<sup>[def. 30.14]</sup> and the *log-loss*<sup>[def. 30.13]</sup> 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<sup>[def. 4.27]</sup> and does not include soft target<sup>[def. 4.28]</sup>.

**CE Loss with Sigmoid Activation**  
**Corollary 30.1**  
**Cross-entropy/Bernoulli loss with Sigmoid Activation:**  
If the last layer uses a sigmoid activation<sup>[def. 29.3]</sup> the log-loss becomes:

$$l_y(\mathbf{z}) = -y \ln \sigma(z) - (1 - y) \ln(1 - \sigma(z)) \tag{30.10}$$

**CE Loss with Softmax Activation**  
**Corollary 30.2** [proof 37.13]  
**Cross-entropy/Log loss with Softmax Activation:**  
The cross entropy loss function is usually used together with softmax-activation functions<sup>[def. 29.6]</sup> in order to measure how close the output probabilities are to the real distribution. In this case the cross-entropy/log-loss becomes:

$$l_y(\mathbf{z}) = -\sum_{k=1}^K y_k \cdot \ln \sigma_k^{\max}(\mathbf{z}) = -\sum_{k=1}^K y_k z_k + \ln \left( \sum_{l=1}^K e^{(z_l)} \right)$$

hard targets  
Def. 4.27  $-z_y + \ln \left( \sum_{l=1}^K e^{(z_l)} \right)$  (30.11)

$$\frac{\partial}{\partial z_k} l_y(\mathbf{z}) = \sigma^{\max}(z_k) - y_k = \hat{p}_k - y_k \tag{30.12}$$

**add connection to sigmoid and motivation: deep learning book ps 182-183**  
**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 30.15** [proof 37.14]  
**Binary Cross Entropy (BCE) Loss:**  
In case of only two classes the cross entropy loss functions can be written as:

$$l_y(\hat{\mathbf{y}}) = -\ln f(y|\hat{\mathbf{y}}; \Theta) = -[y \ln \hat{p}_k + (1 - y) \ln(1 - \hat{p}_k)] \text{ for } y \in \{0, 1\} \tag{30.13}$$

- 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**  
**move to optimisation?**  
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  $\eta$ .
- If the gradient is very/too big  $\Rightarrow$  decrease learning rate  $\eta$ .

**Time Based Decay**  
**Proposition 30.1 Time Based Decay:**

Shrinks the learning rate with time more and more.

$$\eta_t = \eta_t \cdot \frac{1}{(1 + \text{decay} \cdot t)} \tag{30.14}$$

**Proposition 30.2 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:

$$\eta_t = \min(C_1, C_2) = \min\left(C_1, \frac{C_1 t'}{t}\right) \tag{30.15}$$

**Step Based Decay**  
**Proposition 30.3 Step Based Decay:** Divide the learning rate by a fixed factor  $f$  after every  $t'$  epochs:

$$\eta_t = \eta_0 \cdot f^{\left(\left\lfloor \frac{t}{t'} \right\rfloor\right)} \tag{30.16}$$

Exponential Based Decay



Neural Network Layers

- 1. Linear Layers
- 2. Convolutional Layers

Recurrent Neural Networks (RNN)

Goal

Given an observation *sequence* of inputs  $\mathbf{x}^1, \dots, \mathbf{x}^s$  we want to *learn* the transition<sup>[def. 12.5]</sup> (and eventually output<sup>[def. 12.6]</sup> )model of a *discrete time state space model*<sup>[def. 12.7]</sup> that:  

①

Is *Markovian*<sup>[def. 13.1]</sup>  

②

Is *Stationary Time-homogeneous/-invariant* (see also definition 13.15)

$\mathbf{h}^{(t)} = \mathbf{F}(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \theta_{\mathbf{F}})$   
 $t = 1, \dots, s$

$\mathbf{h}^0 = \mathbf{0}$   
 $t = 1, \dots, s$

(32.1)

1. Simple RNNs

Definition 32.1 Recurrent Neural Network (RNN):  
Is a *linear discrete time state space model* <sup>[def. 12.6]</sup> with element-wise non-linearities:  
State Transition Model:  

$\mathbf{F}(\mathbf{h}, \mathbf{x}; \theta_{\mathbf{F}}) = \mathbf{W}\mathbf{h} + \mathbf{U}\mathbf{x} + \mathbf{b}$   
 $\mathbf{F} = \varphi \circ \mathbf{F}$

$\theta_{\mathbf{F}} = (\mathbf{U}, \mathbf{W}, \mathbf{b})$

(32.2)

Output Model:  

$\mathbf{H}(\mathbf{h}; \theta_{\mathbf{H}}) := \mathbf{V}\mathbf{h} + \mathbf{c}$   
 $\mathbf{y} = \bar{\mathbf{H}}(\mathbf{h}; \theta_{\mathbf{H}}) = \varphi_{\mathbf{H}} \circ \mathbf{H}$

$\theta_{\mathbf{H}} = (\mathbf{V}, \mathbf{c})$

(32.3)

The diagram illustrates an unrolled RNN across three time steps:  $t-1$ ,  $t$ , and  $t+1$ . Each time step is represented by a vertical capsule containing four circles, representing the hidden state. The capsules are connected sequentially by horizontal arrows labeled  $\mathbf{W}$ , representing the state transition. The first capsule has a self-loop arrow labeled  $\mathbf{W}$ . Below each capsule, an input vector  $\mathbf{x}$  (blue circle) is shown with an upward arrow labeled  $\mathbf{U}$ . Above each capsule, an output vector  $\hat{\mathbf{y}}$  (green circle) is shown with a downward arrow labeled  $\mathbf{V}$ . The sequence of inputs is  $\mathbf{x}^{(t-1)}, \mathbf{x}^{(t)}, \mathbf{x}^{(t+1)}$  and the sequence of outputs is  $\hat{\mathbf{y}}^{(t-1)}, \hat{\mathbf{y}}^{(t)}, \hat{\mathbf{y}}^{(t+1)}$ . Ellipses ( $\dots$ ) indicate the continuation of the sequence before and after the shown steps.

Main difference to FCFFN<sup>[def. 30.8]</sup>

- The weights/parameters of the different layers are *shared*!
- Inputs  $\mathbf{x}_i$  are processed in a sequence.

1.1. Backpropagation Trough Time (BPTT)

2. Gated Recurrent Units (GRUs)

3. Long Term Short Term Memory (LSTM)



Learning with General Graphs

1. Undirect Graphs

1.1. 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:

$$\mathbf{X}^{l+1} = \varphi\left(\mathbf{W}^l \mathbf{X}^l\right) \quad \text{s.t.} \quad \mathbf{X}^L = \mathbf{Y} \tag{34.1}$$

We now want to model the relationship between the different samples a.k.a. nodes by introducing a matrix  $\mathbf{Q}$ :

$$\boxed{\mathbf{X}^{l+1} = \varphi\left(\mathbf{W}^l \mathbf{X}^l \mathbf{Q}\right)} \tag{34.2}$$

$\mathbf{X}^l \mathbf{Q}$  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 34.1 Linear Shift Invariant Filter for Graphs:** Is a linear function  $\mathbf{H}$  over a graph with adjacency matrix  $\mathbf{A}$  (possibly degree-normalized  $\bar{\mathbf{A}}$ ) s.t.

$$\mathbf{H}(\mathbf{A}\mathbf{x}) = \mathbf{A}(\mathbf{H})\mathbf{x} \iff \mathbf{H} \text{ and } \mathbf{A} \text{ commute} \tag{34.3}$$

**Theorem 34.1 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:

$$\mathbf{H} = \theta_0 \mathbf{I} + \theta_1 \mathbf{A} + \theta_2 \mathbf{A}^2 + \dots + \theta_n \mathbf{A}^n \tag{34.4}$$

1.1.2. Building the Smoothing Matrix

① Favor normalized adjacency matrices:

$$\bar{\mathbf{A}} := \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \quad \text{with} \quad \mathbf{D} \text{ degree matrix}^{\text{[def. 18.7]}} \tag{34.5}$$

② Restrict to linear shift invariant filters of order 1<sup>[def. 34.1]</sup>:

$$\mathbf{H}(\theta_0, \theta_1) = \theta_0 \mathbf{I} + \theta_1 \bar{\mathbf{A}} \tag{34.6}$$

③ Use only one parameter  $\theta_0 = \theta_1$

④ Stack such filters in order to regain high order:

$$\mathbf{H}^k(\theta_0, \theta_1) = \underline{\theta^k \mathbf{I}} + \underline{\theta^k \bar{\mathbf{A}}} \tag{34.7}$$

1.1.3. Problems

If we apply an operator and the operator norm?? is larger than one, then applying the operator repeatedly may lead to blow up.

**Lemma 34.1 Eigenvalues of Graph Convolution Layers:** Let  $\lambda_1 \geq \dots \geq \lambda_n$  be the eigenvalues of  $\bar{\mathbf{A}}$  then it holds that:

$$1 = \lambda_1 \geq \lambda_n \geq -1 \quad \lambda_1, \dots, \lambda_n \in \text{eigenval}(\bar{\mathbf{A}}) \tag{34.8}$$

**Corollary 34.1 :**

$$\underline{\mathbf{I} + \bar{\mathbf{A}} \in [0, 2]} \implies \text{possibly blow up} \tag{34.9}$$

1.1.4. Semi-Supervised Classification with GCNNs

Solution to Corollary 1.1.3

Define a degree normalized matrix by using self-loops:

$$\bar{\mathbf{A}} + \mathbf{I} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} + \mathbf{I} \xrightarrow{\text{instead}} \widetilde{\mathbf{D}}^{-\frac{1}{2}} \widetilde{\mathbf{A}} \widetilde{\mathbf{D}}^{-\frac{1}{2}} =: \mathbf{Q}$$

**Definition 34.2 Self loop Adjacency Matrix:**

$$\bar{\mathbf{A}} := \mathbf{A} + \mathbf{I}_n \tag{34.10}$$

**Definition 34.3 Self loop Degree Matrix:**

$$\widetilde{\mathbf{D}} := \mathbf{D} + \mathbf{I}_n \iff \tilde{d}_i = 1 + \sum_{\substack{j < i \\ j}} \mathbf{A}_{ij} \tag{34.11}$$

**Definition 34.4 Coupling Matrix**

**Q:**

$$\mathbf{Q} = \widetilde{\mathbf{D}}^{-\frac{1}{2}} \widetilde{\mathbf{A}} \widetilde{\mathbf{D}}^{-\frac{1}{2}} \tag{34.12}$$
  
$$q_{ij} = \frac{a_{ij} + \delta_{ij}}{\sqrt{\tilde{d}_i \tilde{d}_j}} \quad \text{wit} \quad \delta_{ij} = \begin{cases} 1 & \text{if } j = i \\ 0 & \text{else} \end{cases}$$

**Explanation 34.1** (Coupling Matrix).

$$\left(\mathbf{x}^l \mathbf{Q}\right)_{ij} = \left(\begin{bmatrix} \mathbf{x}_1^l & \dots & \mathbf{x}_n^l \end{bmatrix} \mathbf{Q}\right)_{ij} = \sum_{k=1}^n x_{ik}^l q_{kj} \tag{34.13}$$

$$\left(\mathbf{x}^l \mathbf{Q}\right)_{:j} = \sum_{k=1}^n q_{kj} \cdot \mathbf{x}_k^l \tag{34.14}$$

Thus the coupling matrix<sup>[def. 34.4]</sup> leads to a smoothed combination with the neighboring nodes.

# Regularization

## 1. Weight Decay

### 1.1. L1/Lasso Regularization

### 1.2. L2 Regularization

## 2. 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 35.1 Data Augmentation:** Is augmentation of the training data by artificially generated training data.

**Explanation 35.1.** *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.

### 2.1. Generation By Invariant Transformations

**Formula 35.1** [example 35.2]

**Invariance Augmentation:**

① find transformations  $\tau$  that our original data should be invariant too.

② Generate *virtual examples* by applying  $\tau$  to our original training set:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n \xrightarrow{\tau} \{\tau(\mathbf{x}_i), \mathbf{y}_i\}_{i=1}^n \quad (35.1)$$

## 3. Task Augmentation

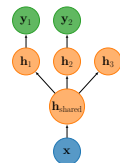
**Definition 35.2** [example 35.3]

**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.

## 4. Multi Task Learning

**Definition 35.3 Multi Task Learning:**

Share representation of the sub-model across different tasks and learn the model jointly i.e. by a combined objective.



## 5. Examples

**Example 35.1 2D Invariances:**

- Scale Changes
- Rotations and reflections
- Transformations s.a. PCA

**Example 35.2 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).

**Example 35.3 Pre-trained Models:**

- Pre-trained Word Embeddings
- Pre-trained filters and maps

## 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 36.1 Memory Networks:** Memory components and a model that is trained to effectively operate with that memory component.



# Proofs

## Activation Functions

Definition 29.4

$$\sigma'(z) = \frac{\partial}{\partial z} \left(1 + e^{-z}\right)^{-1} = -\left(1 + e^{-z}\right)^{-2} \frac{\partial}{\partial z} \left(1 + e^{-z}\right)$$

$$= -\frac{1}{(1 + e^{-z})^2} e^{-z} (-1) = \frac{1}{(1 + e^{-z})^2} \left(1 - e^{-z}\right)$$

$$= \frac{1}{(1 + e^{-z})} - \frac{1}{(1 + e^{-z})^2} = \sigma(z) - \sigma(z)^2$$

$$\sigma'(z) = \frac{\partial}{\partial z} \left(1 + e^{-z}\right)^{-1} = -\left(1 + e^{-z}\right)^{-2} e^{-z} (-1)$$

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

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

Proof 37.2:
Property 29.1

$$\sigma\left(\ln\left(\frac{t}{1-t}\right)\right) = \frac{1}{\frac{t}{1-t}} = t$$
(37.1)

Proof 37.3:
Property 29.2

$$2\sigma(2z) - 1 = 2\frac{1}{1 + e^{-2z}} - 1 = \frac{2}{1 + e^{-2z}} - \frac{1 + e^{-2z}}{1 + e^{-2z}}$$

$$= \frac{1 - e^{-2z}}{1 + e^{-2z}} = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \tanh(z)$$

Proof 37.4:
Definition 29.5

$$\tanh(z) = \frac{\sinh(z)}{\cosh(z)} = \frac{\frac{\partial}{\partial z} \sinh(z) \cdot \cosh(z) - \frac{\partial}{\partial z} \cosh(z) \cdot \sinh(z)}{\cosh^2(z)}$$

$$= \frac{\cosh^2(z) - \sinh^2(z)}{\cosh^2(z)} = 1 - \frac{\sinh^2(z)}{\cosh^2(z)}$$

Proof 37.5:
Definition 29.6

Softmax function derivative
 $\nabla_{z_j} \sigma_i^{\max}(z)$ :

$$i \neq j :$$

$$\nabla_{z_j} \frac{\exp(\mathbf{x}^\top \theta_i)}{\sum_{k=1}^c \exp(\mathbf{x}^\top \theta_k)} =$$
Q.R.
$$= \frac{0 - e^{x_j} e^{x_i}}{\left(\sum_{k=1}^c e^{x_k}\right)^2} = -\frac{e^{x_j}}{\sum_{k=1}^c e^{x_k}} \cdot \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}}$$

$$= -\sigma_i^{\max}(z) \sigma_j^{\max}(z)$$

$$i = j :$$

$$\nabla_{z_i} \frac{\exp(\mathbf{x}^\top \theta_i)}{\sum_{k=1}^c \exp(\mathbf{x}^\top \theta_k)} =$$

$$= \frac{e^{x_i} \sum_{k=1}^c e^{x_k} - e^{x_i} e^{x_i}}{\left(\sum_{k=1}^c e^{x_k}\right)^2}$$

$$= \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} - \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}} \cdot \frac{e^{x_i}}{\sum_{k=1}^c e^{x_k}}$$

$$= \sigma_i^{\max}(z) - \left(\sigma_i^{\max}(z)\right)^2$$

Proof 37.6:
Corollary 37.6

$$\sigma_1^{\max}(\mathbf{x}) = \frac{e^{\mathbf{x}^\top \theta_1}}{e^{\mathbf{x}^\top \theta_1} + e^{\mathbf{x}^\top \theta_2}} = \frac{e^{-\mathbf{x}^\top \theta_1}}{e^{-\mathbf{x}^\top \theta_1}} \frac{e^{\mathbf{x}^\top \theta_1}}{e^{\mathbf{x}^\top \theta_1} + e^{\mathbf{x}^\top \theta_2}}$$

$$= \frac{e^{\mathbf{x}^\top (\theta_1 - \theta_1)}}{e^{\mathbf{x}^\top (\theta_1 - \theta_1)} + e^{\mathbf{x}^\top (\theta_2 - \theta_1)}} = \frac{1}{1 + e^{-\mathbf{x}^\top \theta}}$$

Proof 37.7:
Definition 29.7

$$y = \frac{e^x}{e^x} \cdot \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x} = \frac{1 + e^x - 1}{1 + e^x} = 1 - \frac{1}{1 + e^x}$$

$$\Rightarrow \frac{1}{1 + e^x} = 1 - y$$

$$1 + e^x = \frac{1}{1 - y} = \frac{1 - y + y}{1 - y} = 1 + \frac{y}{1 - y}$$

$$e^x = \frac{y}{1 - y} \quad \Rightarrow \quad x = \ln \frac{y}{1 - y}$$

Proof 37.8
Stable Softmax<sup>[def. 29.8]</sup>:
We simply multiply the numerator and denominator by a constant  $\hat{a}$ :
$$P_i = \frac{\hat{a} e^{z_i}}{\hat{a} \sum_{k=1}^c e^{z_k}} = \frac{e^{z_i + \log(\hat{a})}}{\sum_{k=1}^c e^{z_k + \log(\hat{a})}}$$
And now we define  $a := \log \hat{a} = -\max(\mathbf{z})$

Proof 37.9:
Corollary 29.10

If  $z \geq 0 \quad \Rightarrow \quad (z)_+ = 2z - z = z$ 
if  $z < 0 \quad \Rightarrow \quad (z)_+ = 0$ 
(37.2)
(37.3)

## Losses

Proof 37.10:
Definition 30.13

$$l_y(\hat{\mathbf{y}}) = -\log f(y|\hat{\mathbf{y}}; \Theta) = -\log \prod_{i=1}^k \hat{y}_i^{\delta[y=c_i]}$$

$$= -\sum_{i=1}^k \delta_{[y=c_i]} \cdot \log \hat{p}_i = \delta_{[y=c_i]} \cdot \log \hat{p}_i$$

Proof 37.11:
Defintion 30.2

$$l_y(\mathbf{z}) = -\sum_{k=1}^K y_k \cdot \ln \sigma_k^{\max}(\mathbf{z}) = -\sum_{k=1}^K y_k \left(z_k - \ln \sum_{l=1}^K e^{z_l}\right)$$

$$= -\sum_{k=1}^K y_k z_k + \sum_{k=1}^K y_k \ln \left(\sum_{l=1}^K e^{z_l}\right)$$

$$= -\sum_{k=1}^K y_k z_k + \ln \left(\sum_{l=1}^K e^{z_l}\right) \cdot \underbrace{\sum_{k=1}^K y_k}_{=1}$$

$$= -\sum_{k=1}^K y_k z_k + \ln \left(\sum_{l=1}^K e^{z_l}\right)$$

Proof 37.12
CE Derivative<sup>[def. 30.14]</sup>:

$$\frac{\partial}{\partial z_k} \text{CE} = \frac{\partial}{\partial z_k} - \sum_{k=1}^K p_k \ln \hat{p}(z_k)$$
(37.4)
$$= -\sum_{k=1}^K p_k \frac{1}{\hat{p}(z_k)} \frac{\partial}{\partial z_k} \hat{p}(z_k)$$
(37.5)

Proof 37.13
Deriv. of CE with Softmax<sup>[cor. 30.2]</sup> and ?? 37.11:

$$\frac{\partial}{\partial z_k} l_y(\mathbf{z}) \stackrel{\text{eq. (30.8)}}{=} \sum_{j=1}^K y_j \frac{1}{\sigma^{\max}(z_j)} \frac{\partial}{\partial z_k} \sigma^{\max}(z_j)$$

$$\stackrel{\text{eq. (29.8)}}{=} -y_k \frac{\sigma^{\max}(z_k) (1 - \sigma^{\max}(z_k))}{\sigma^{\max}(z_k)} \quad j = k$$

$$- \sum_{j \neq k}^K y_j \frac{-\sigma^{\max}(z_k) \sigma^{\max}(z_j)}{\sigma^{\max}(z_j)} \quad j \neq k$$

$$= -y_k (1 - \sigma^{\max}(z_k)) + \sum_{j \neq k} y_j \sigma^{\max}(z_k)$$

$$= -y_k + \sigma^{\max}(z_k) y_k + \sum_{j \neq k} y_j \sigma^{\max}(z_k)$$

$$= \sigma^{\max}(z_k) - y_k = \hat{p}_k - y_k$$

Proof 37.14
Binary Cross Entropy Loss<sup>[def. 30.15]</sup>:
The binary cross entropy can be derived directly from the Bernoulli Distribution??, where y signifies that we penalize only the true label:
$$L_y = \hat{p}^y \cdot (1 - \hat{p})^{1-y} \quad \text{for } y \in \{0, 1\}$$

$$l_y = -\log L_y = -[y \ln \hat{p} + (1 - y) \ln(1 - \hat{p})]$$

## Autoencoders

Proof 37.15
[Corollary 27.1]:
Let  $\mathbf{C}^*$ ,  $\mathbf{D}^*$  be the global solution of problem Equation (27.2). Then we may define a matrix  $\mathbf{A}$  s.t.:
$$\mathbf{D}^* \mathbf{C}^* = \mathbf{D}^* (\mathbf{A}^{-1} \mathbf{A}) \mathbf{C}^* = \widehat{\mathbf{D}}^* \tilde{\mathbf{C}}^* \quad (37.6)$$

Proof 37.16
Proposition 27.2:
$$\mathbf{DCX} \stackrel{??}{=} \mathbf{DCU} \Sigma \mathbf{V}^H = \mathbf{U}^\top \underbrace{\mathbf{U}_k \mathbf{U}}_{\stackrel{??}{=} \mathbf{I}_m} \Sigma \mathbf{V}^H$$

$$= \mathbf{U}_k \Sigma_m \mathbf{V}_m^H = \mathbf{X}_m$$

Examples

Losses

**Example 38.1 Categorical Loss**<sup>[def. 30.13]</sup>:  
Given an observation  $\mathbf{y} = [0 \ 1 \ 0 \ 0 \ 0]^\top$  the *categorical loss*<sup>[def. 30.13]</sup> 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  $\alpha_{1,3,4,5} \in [0, 1]$ .