

Natural Language Processing with Deep Learning

Sentiment Analysis with Machine Learning



Navid Rekab-Saz

navid.rekabsaz@jku.at

Agenda

- Introduction to Machine Learning
- Sentiment Analysis
- Feature Extraction
- *Breaking* the curse of dimensionality!

Agenda

- **Introduction to Machine Learning**
- Sentiment Analysis
- Feature Extraction
- *Breaking* the curse of dimensionality!

Notation

- $a \rightarrow$ a value or a scalar
- $\mathbf{b} \rightarrow$ an array or a vector
 - i^{th} element of \mathbf{b} is the scalar b_i
- $\mathcal{C} \rightarrow$ a set of arrays or a matrix
 - i^{th} vector of \mathcal{C} is \mathbf{c}_i
 - j^{th} element of the i^{th} vector of \mathcal{C} is the scalar $c_{i,j}$

Linear Algebra – Recap

- Transpose

- \mathbf{a} is in $1 \times d$ dimensions $\rightarrow \mathbf{a}^T$ is in $d \times 1$ dimensions
- \mathbf{A} is in $e \times d$ dimensions $\rightarrow \mathbf{A}^T$ is in $d \times e$ dimensions

- Inverse of the square matrix \mathbf{S} is \mathbf{S}^{-1}

- Dot product

- $\mathbf{a} \cdot \mathbf{b}^T = c$
dimensions: $1 \times d \cdot d \times 1 = 1$
- $\mathbf{a} \cdot \mathbf{B} = \mathbf{c}$
dimensions: $1 \times d \cdot d \times e = 1 \times e$
- $\mathbf{A} \cdot \mathbf{B} = \mathbf{C}$
dimensions: $l \times m \cdot m \times n = l \times n$

Statistical Learning

- Given N observed data points:

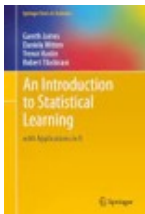
$$\mathbf{X} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}]$$

accompanied with output (label) values:

$$\mathbf{y} = [y_1, y_2, \dots, y_N]$$

and each data point is defined as a vector with l dimensions (features):

$$\mathbf{x}^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_l^{(i)}]$$



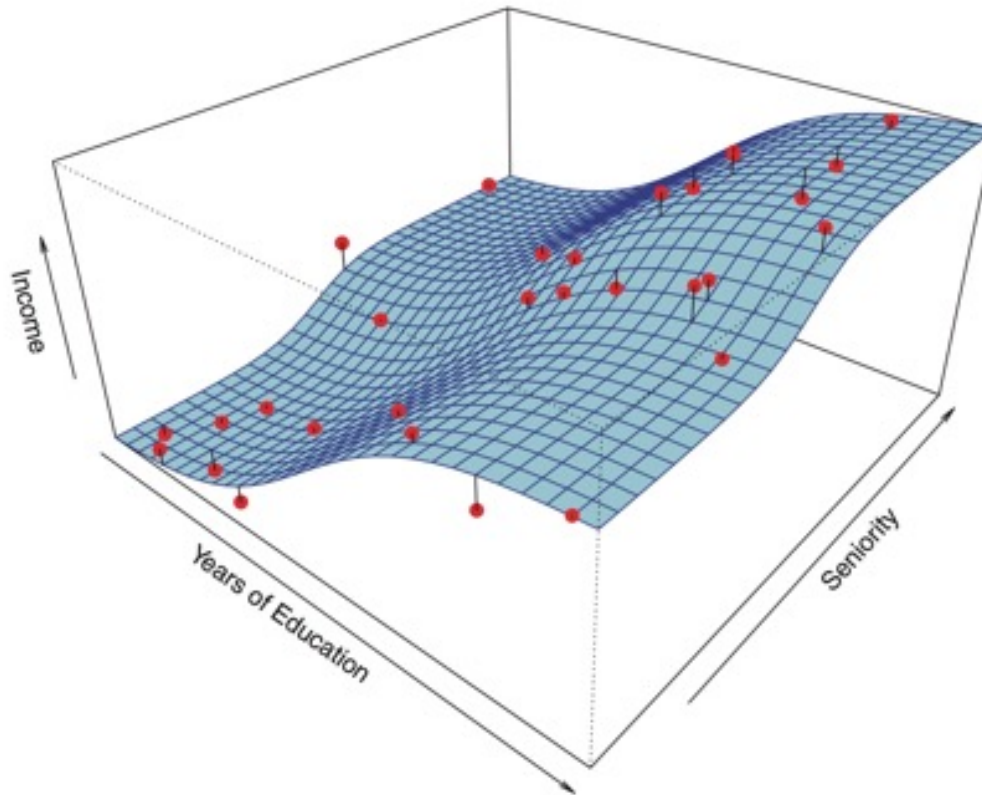
Statistical Learning

- Statistical learning assumes that there exists a **TRUE function** (f^{TRUE}) that has generated these data:

$$\mathbf{y} = f^{TRUE}(\mathbf{X}) + \epsilon$$

- f^{TRUE}
 - The *true* but *unknown* function that produces the data
 - A fixed function
- $\epsilon > 0$
 - Called **irreducible error**
 - Rooted in the constraints in gathering data, and measuring and quantifying features

Example f^{TRUE}



$f^{TRUE} \rightarrow$ blue surface

$X \rightarrow$ Red points with two features: **Seniority**, **Years of Education**

$y \rightarrow$ Income

$\epsilon \rightarrow$ the differences between the data points and the surface

Machine Learning Model

- A machine learning (ML) model tries to estimate f^{TRUE} by defining function f :

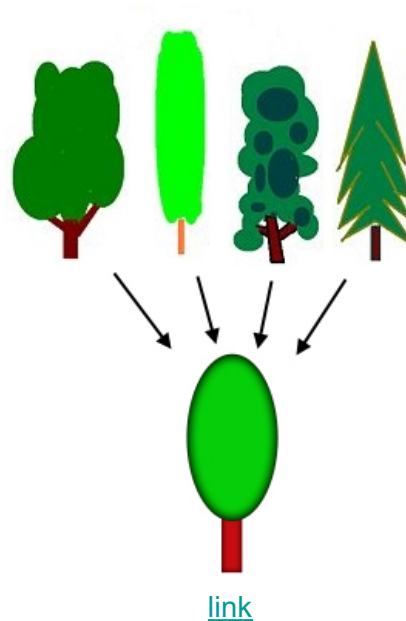
$$\hat{\mathbf{y}} = f(\mathbf{X})$$

such that $\hat{\mathbf{y}}$ (predicted outputs) be close to \mathbf{y} (real outputs).

- The differences between the values of $\hat{\mathbf{y}}$ and \mathbf{y} is **reducible error**
 - Can be reduced by better models, better estimations of f^{TRUE}

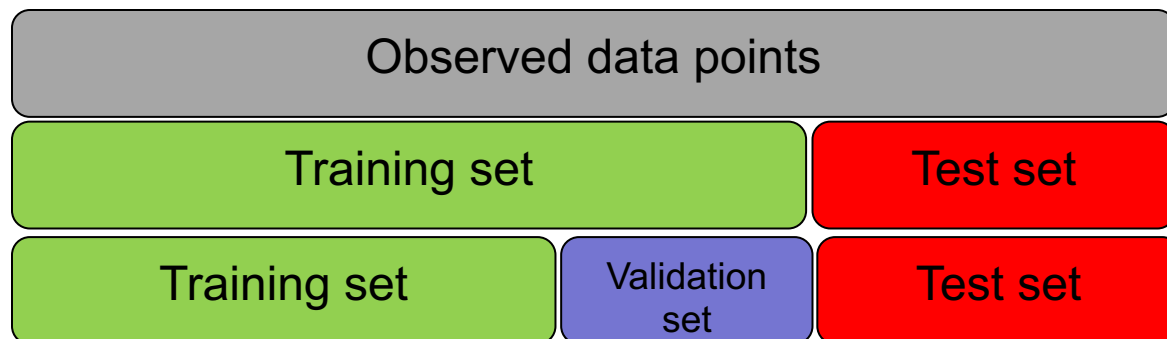
Generalization

- The aim of machine learning is to create a model using **observed experiences** (training data) that generalizes to the problem domain, namely **performs** well on **unobserved instances** (test data)



Learning the model – Splitting dataset

- Data points are splitted into:
 - **Training set**: for training the model
 - **Validation set**: for tuning model's hyper-parameters
 - **Test set**: for evaluating model's performance
- Common train – validation – test splitting sizes
 - 60%, 20%, 20%
 - 70%, 15%, 15%
 - 80%, 10%, 10%



Learning the model

Dataset

Features / Variables (X)				Labels / Output Variable (Y)
sex	age	Pstatus	romantic	Walc
F	18	A	no	1
F	17	T	no	1
F	15	T	no	3
F	15	T	yes	1
F	16	T	no	2
M	16	T	no	2
M	16	T	no	1
F	17	A	no	1
M	15	A	no	1
M	15	T	no	1
F	15	T	no	2
F	15	T	no	1
M	15	T	no	3
M	15	T	no	2
M	15	A	yes	1
F	16	T	no	2
F	16	T	no	2
F	16	T	no	1
M	17	T	no	4

Pstatus: parent's cohabitation status ('T' - living together 'A' - apart)
Romantic: with a romantic relationship
Walc: weekend alcohol consumption (from 1 - very low to 5 - very high)

<http://archive.ics.uci.edu/ml/datasets/STUDENT+ALCOHOL+CONSUMPTION#>

Learning the model

Train Set

sex	age	Pstatus	romantic	Walc
F	18	A	no	1
F	17	T	no	1
F	15	T	no	3
F	15	T	yes	1
F	16	T	no	2
M	16	T	no	2
M	16	T	no	1
F	17	A	no	1
M	15	A	no	1
M	15	T	no	1
F	15	T	no	2
F	15	T	no	1
M	15	T	no	3

Test Set

M	15	T	no	2
M	15	A	yes	1
F	16	T	no	2
F	16	T	no	2
F	16	T	no	1
M	17	T	no	4

Learning the model

Train Set

sex	age	Pstatus	romantic	Walc
F	18	A	no	1
F	17	T	no	1
F	15	T	no	3
F	15	T	yes	1
F	16	T	no	2
M	16	T	no	2
M	16	T	no	1
F	17	A	no	1
M	15	A	no	1
M	15	T	no	1
F	15	T	no	2
F	15	T	no	1
M	15	T	no	3

Test Set

M	15	T	no	?
M	15	A	yes	?
F	16	T	no	?
F	16	T	no	?
F	16	T	no	?
M	17	T	no	?

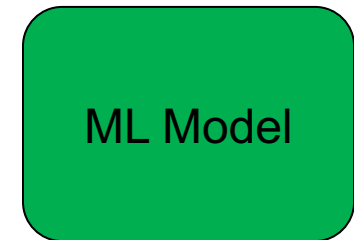
2
1
2
2
1
4

y

Learning the model

Train Set

sex	age	Pstatus	romantic	Walc
F	18	A	no	1
F	17	T	no	1
F	15	T	no	3
F	15	T	yes	1
F	16	T	no	2
M	16	T	no	2
M	16	T	no	1
F	17	A	no	1
M	15	A	no	1
M	15	T	no	1
F	15	T	no	2
F	15	T	no	1
M	15	T	no	3



Test Set

M	15	T	no	?
M	15	A	yes	?
F	16	T	no	?
F	16	T	no	?
F	16	T	no	?
M	17	T	no	?

2
1
2
2
1
4

y

Learning the model

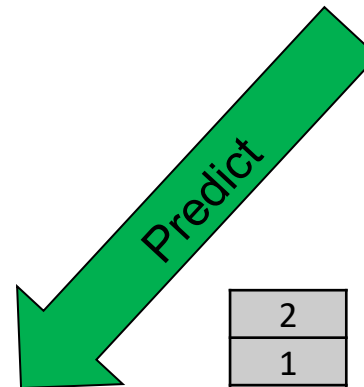
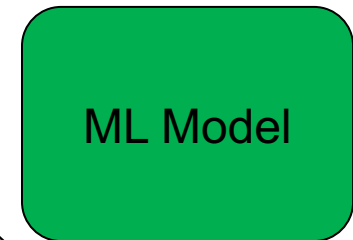
Train Set

sex	age	Pstatus	romantic	Walc
F	18	A	no	1
F	17	T	no	1
F	15	T	no	3
F	15	T	yes	1
F	16	T	no	2
M	16	T	no	2
M	16	T	no	1
F	17	A	no	1
M	15	A	no	1
M	15	T	no	1
F	15	T	no	2
F	15	T	no	1
M	15	T	no	3

Test Set

M	15	T	no	1
M	15	A	yes	1
F	16	T	no	2
F	16	T	no	2
F	16	T	no	3
M	17	T	no	4

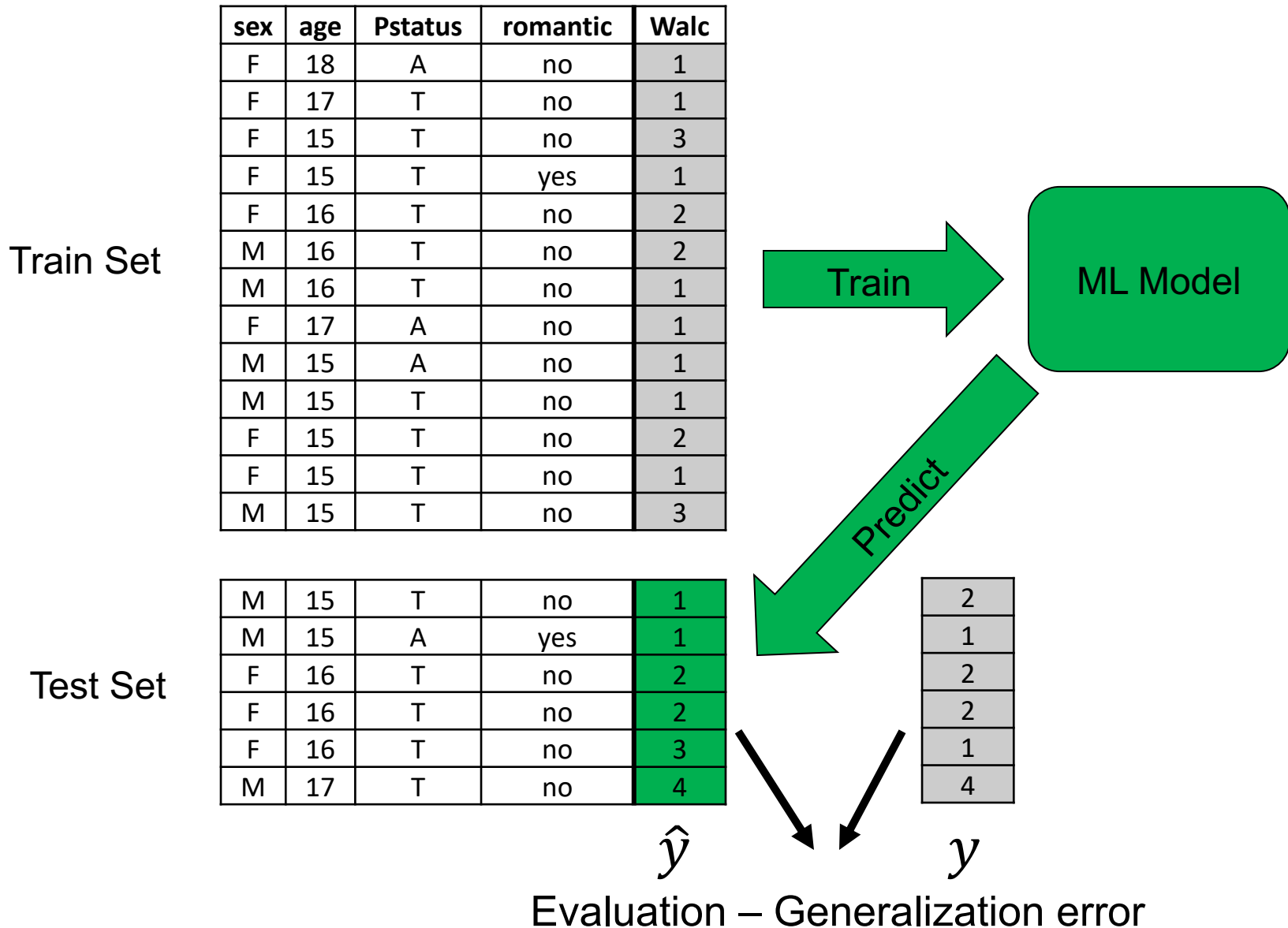
\hat{y}



2
1
2
2
1
4

y

Learning the model



Tuning hyper parameters – Model selection

- Decide on the exploration of several sets of the model's hyper-parameters
- Train a separate model per each set using training set
- Among the trained models, select the best performing one based on the evaluation result on validation set
- Take the selected model and evaluate it on test set → final model performance

ML models

■ Parametric models

- The model is defined as a function (or a family of functions) consisting of a set of parameters
- Functions such as linear regression, logistic regression, naïve Bayes, and neural networks
- The problem of finding the ML model is reduced to finding the optimum values for the parameters

■ Non-parametric models

- There is no assumption about the form of the function
- The model is directly learned from data
- ML models such as SVM, k-NN, smoothing spline, gaussian processes

Term of the day!

Inductive bias: all assumptions we consider in defining and creating an ML model. Our prior knowledge about what f^{TRUE} should be.

A sample ML model: Linear Regression

- f is defined as a Linear Regression function:

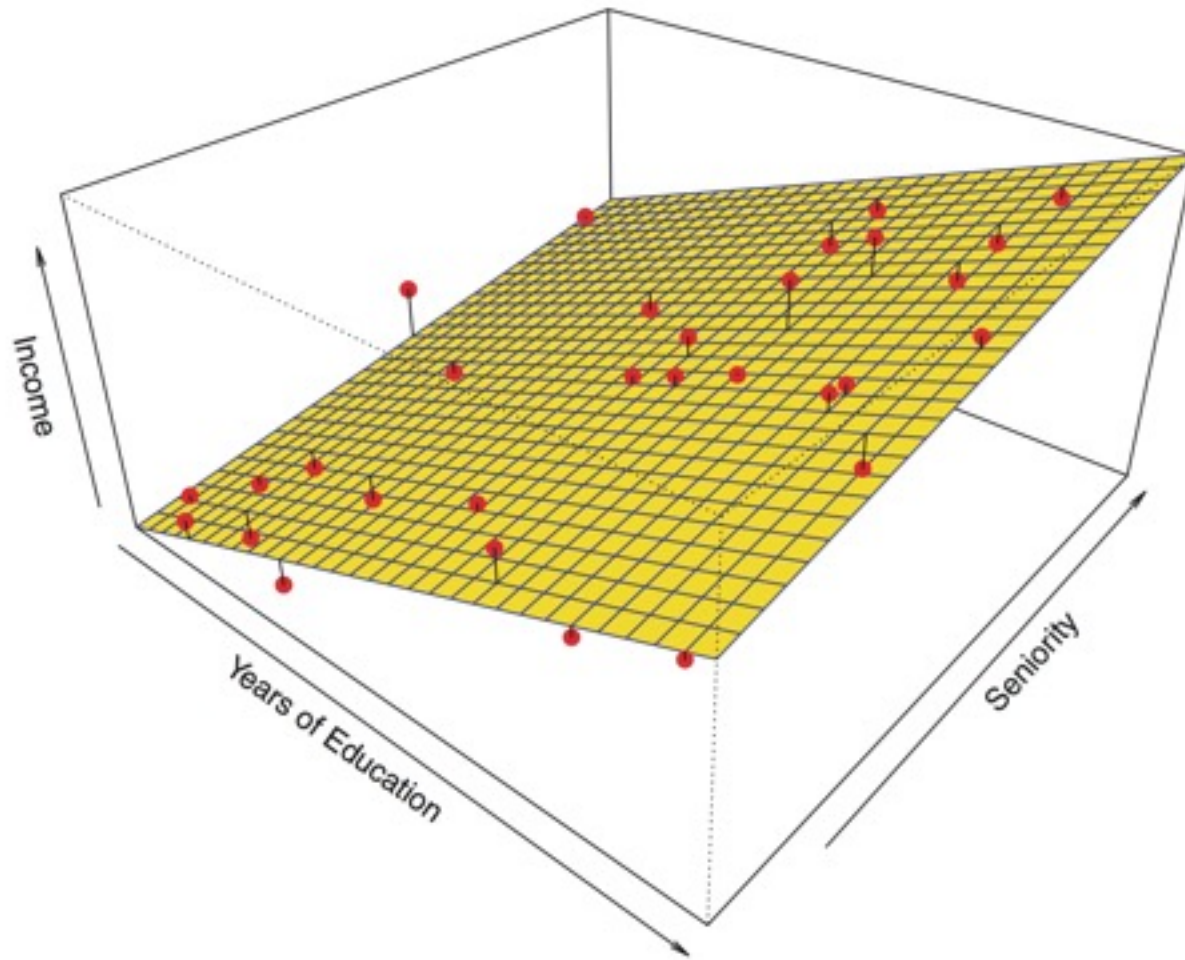
$$y = f(\mathbf{x}; \mathbf{w}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_lx_l$$

where $\mathbf{w} = [w_0, w_1, \dots, w_l]$ is the set of model parameters

- In the “income” example:

$$income = f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 \times education + w_2 \times seniority$$

A trained Linear Regression model



Loss Function

- Optimization of parameters is done by first defining a loss function
- A loss function measures the discrepancies between the predicted outputs $\hat{\mathbf{y}}$ and real ones \mathbf{y}
- E.g. **Mean Square Error (MSE)** – a common regression loss function:

$$\mathcal{L}(y_i, \hat{y}_i; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Loss functions for classification: Next lectures

Good to know! What is Mean Absolute Error and how is it different from MSE?

Optimization

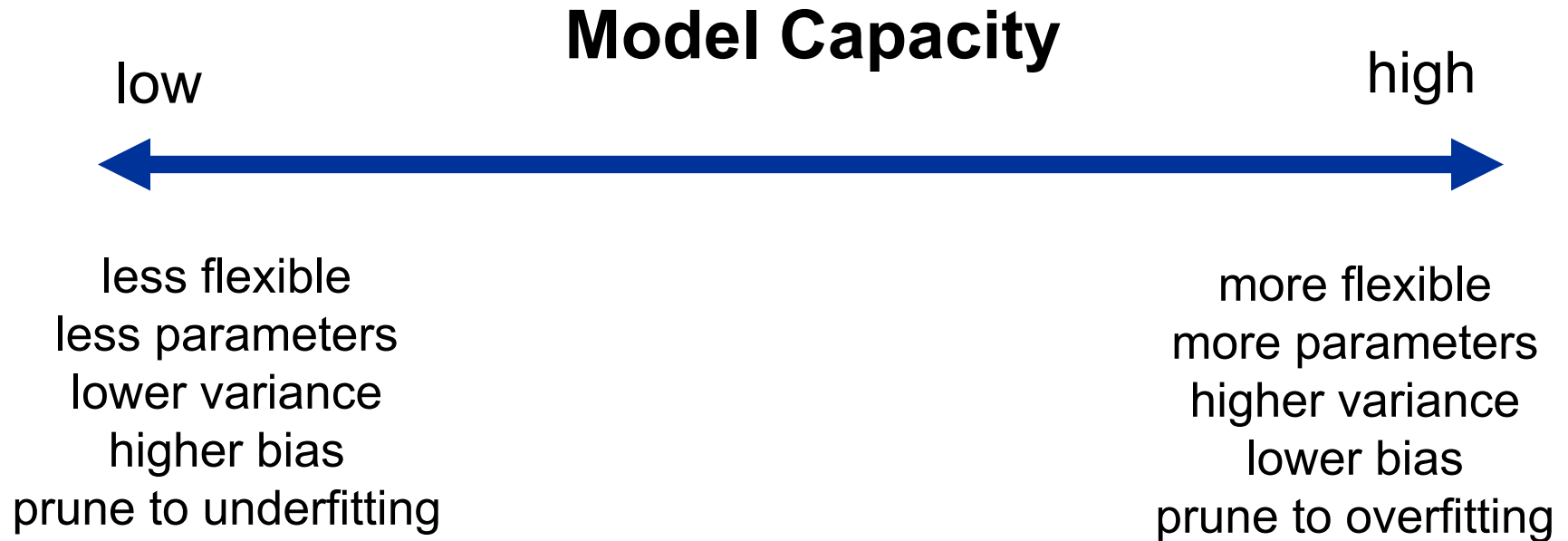
- Next, training data is used to find an *optimum* set of parameters \mathbf{w}^* by optimizing the loss function:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \mathcal{L}(y_i, \hat{y}_i; \mathbf{w})$$

$$\text{MSE: } \mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i; \mathbf{w}))^2$$

- How to optimize:
 - **Stochastically**, e.g. using Stochastic Gradient Descent (SGD)
→ next lecture
 - **Analytically**, e.g. in linear regression → Deep Learning book 5.1.4

ML models... cont.



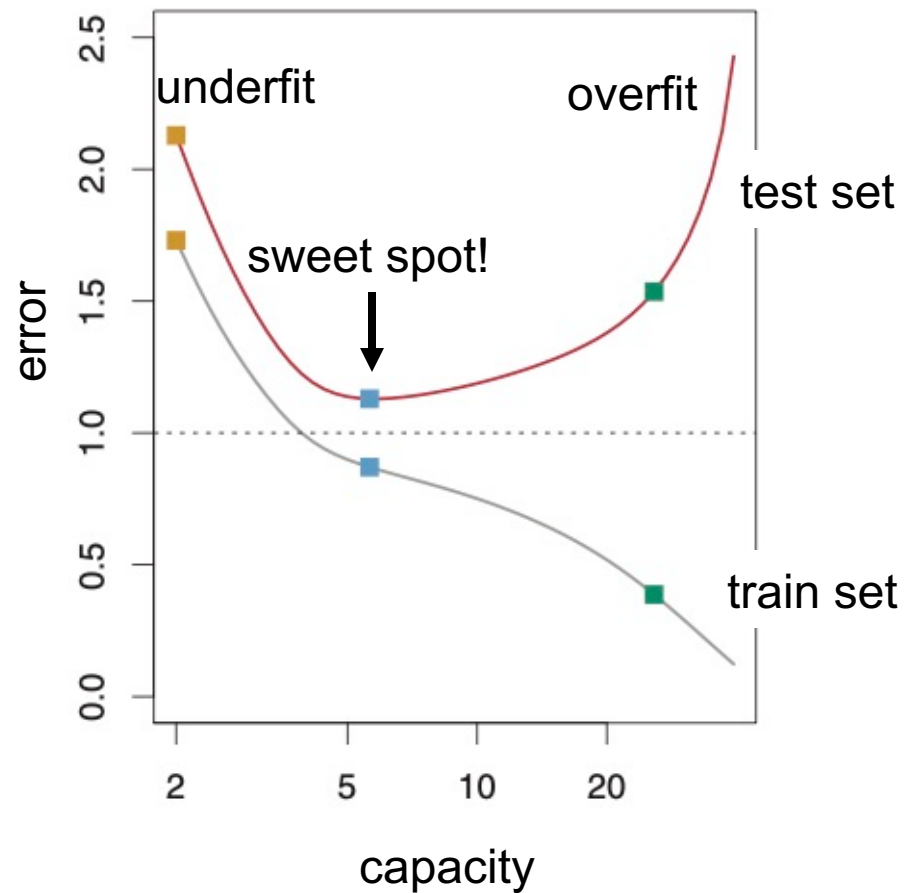
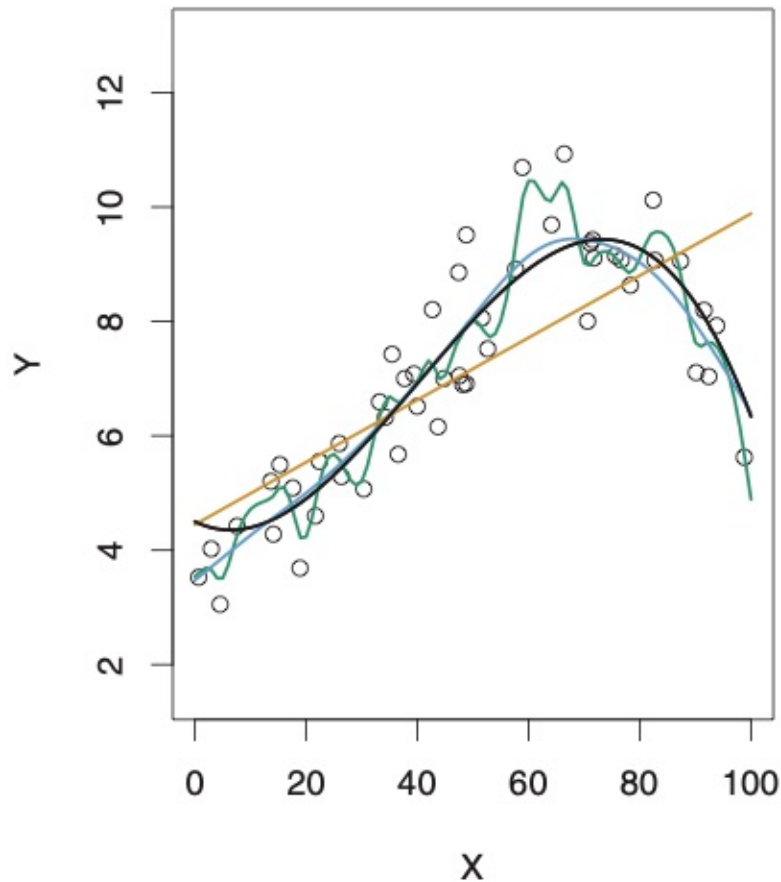
Terms of the day!

(Statistical) Bias indicates the amount of assumptions, taken to define a model. Higher bias means more assumptions and less flexibility, as in linear regression.

Variance: in what extent the estimated parameters of a model vary when the values of data points change (are resampled).

Overfitting: When the model exactly fits to training data, namely when it also captures the noise in data.

Learning Curve



Models:

black $\rightarrow f^{TRUE}$

orange \rightarrow linear regression

blue and green \rightarrow two smoothing spline models

Regularization

- A regularization method introduces additional information (assumptions) to **avoid overfitting** by **decreasing variance**
- E.g. adding the squared L2 norm of parameters to loss function:

$$\mathcal{L}(y_i, \hat{y}_i; \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \|\mathbf{w}\|_2^2$$

$$\|\mathbf{w}\|_2 = \sqrt{\sum_i w_i^2}$$

Common Evaluation Metrics

- Classification

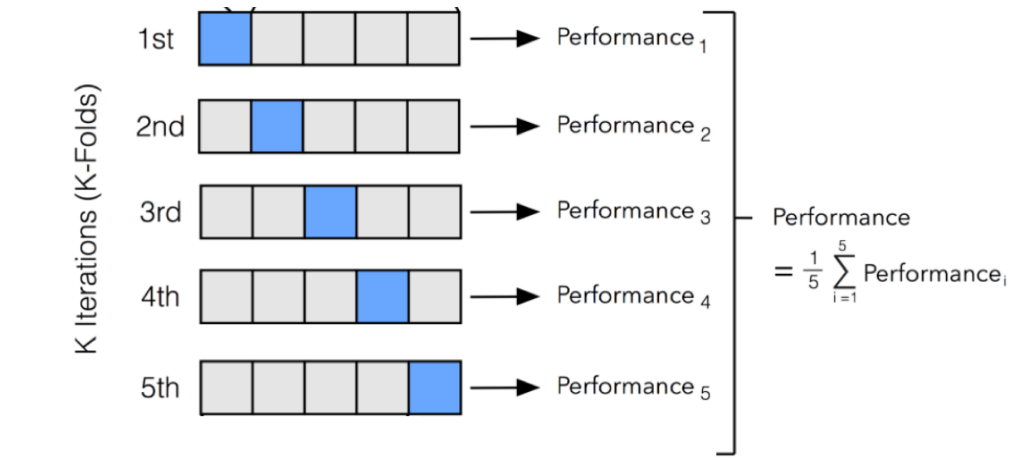
- Accuracy $\frac{\text{\# of correct predictions}}{\text{\# of samples}}$
- Precision $\frac{TP}{TP+FP}$
- Recall $\frac{TP}{TP+FN}$
- F-measure $\frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$

- Regression

- MSE
- R-squared

k-fold Cross Validation

- A rigours evaluation method
 - avoids bias in train/test splitting
- How to
 - Split data into k equal-size folds ($k=5$ or 10)
 - Repeat k times:
 - Use one **left-out** fold for test
 - Use the rest of $k-1$ folds for training
 - Final performance is the **average** of the evaluation results of the k models



Agenda

- Introduction to Machine Learning
- **Sentiment Analysis**
- Feature Extraction
- *Breaking* the curse of dimensionality!

A tough Example!

“This past Saturday, I bought a Nokia phone and my girlfriend bought a Motorola phone with Bluetooth. We called each other when we got home. The voice on my phone was clear, better than my previous Samsung phone. The battery life was however short. My girlfriend was quite happy with her phone. I wanted a phone with good sound quality. So my purchase was a real disappointment. I returned the phone yesterday.”

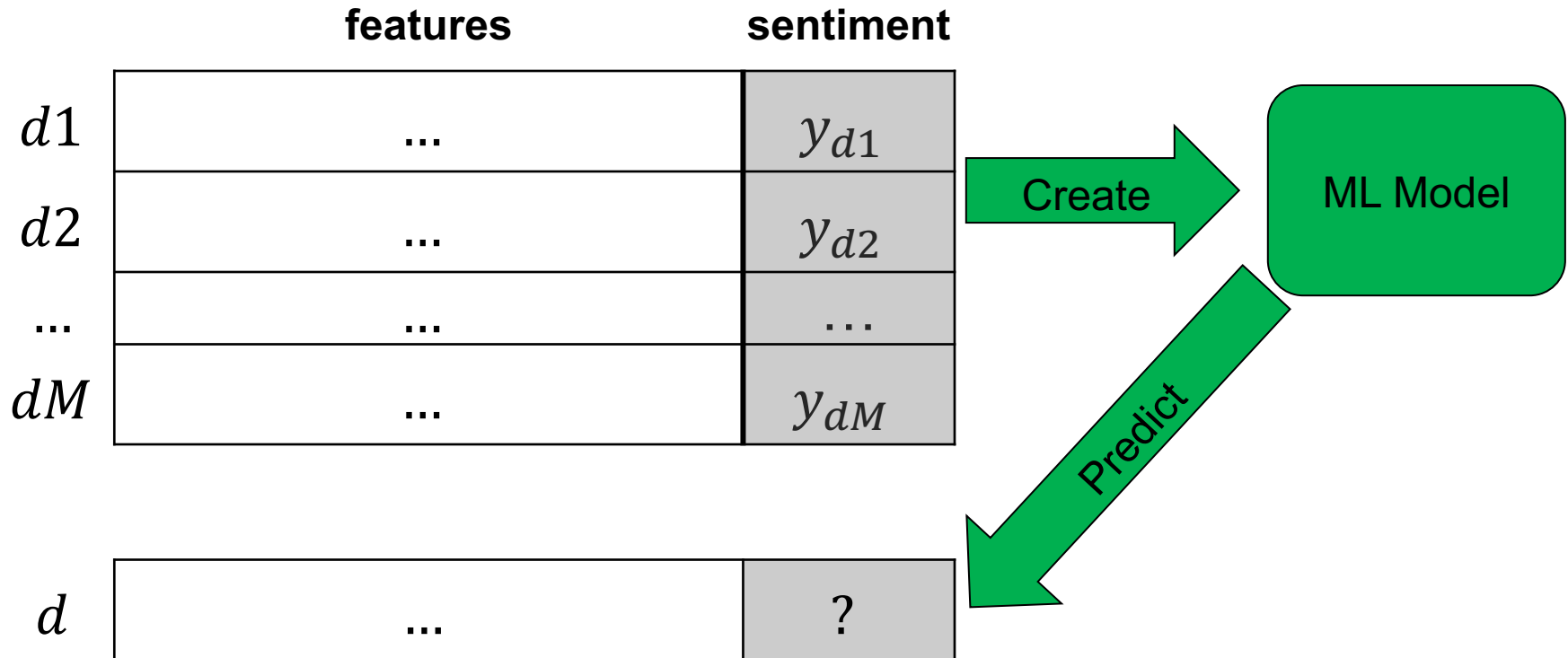
Text-Level Sentiment Analysis

- Text- or document-level sentiment analysis assumes that whole the text expresses **one sentiment** about **one opinion target**
 - Not like the previous example!
- We approach sentiment prediction with ML

Problem definition

- A dataset consist of M text documents and their sentiments (outputs)
- Possible sentiment values:
 - $[-1, 0, 1] \rightarrow$ [negative, neutral, positive] (classification problem)
 - Real-valued numbers e.g. stock price (regression problem)

Sentiment Analysis with ML



Agenda

- Introduction to Machine Learning
- Sentiment Analysis
- **Feature Extraction**
- *Breaking* the curse of dimensionality!

Dictionary

- To extract features, first a dictionary with N words (terms) is defined:

$$[t_1, t_2, \dots, t_N]$$

Document-Term Matrix

- The features are based on the terms in the dictionary
 - Bag of Words (BoW) representations of documents
- $x_{t,d}$ is feature value \rightarrow **weight** of term t in document d

	$t1$	$t2$...	tN	sentiment
$d1$	$x_{t1,d1}$	$x_{t2,d1}$...	$x_{tN,d1}$	y_{d1}
$d2$	$x_{t1,d2}$	$x_{t2,d2}$...	$x_{tN,d2}$	y_{d2}
...
dM	$x_{t1,dM}$	$x_{t2,dM}$...	$x_{tN,dM}$	y_{dM}

Term Weightings

- A term weighting method measures the importance of a term in a document
- One common method is to count the number of occurrences of a term in a document \Rightarrow **term count**:

$$x_{t,d} = tc_{t,d} = \# \text{ of occurrences of } t \text{ in } d$$

- Using **logarithm** to dampening raw counts is shown to be more effective \Rightarrow **term frequency**:

$$x_{t,d} = tf_{t,d} = \log(1 + tc_{t,d})$$

On informativeness of less frequent terms

- Terms that do not appear often usually carry more information in comparison with highly frequent ones
 - e.g., **JKU** in a large news corpora
- Inverse document frequency (idf) is a well-known method to measure how often words appear in a collection:

$$\text{idf}_t = \log\left(\frac{M}{\text{df}_t + 1}\right)$$

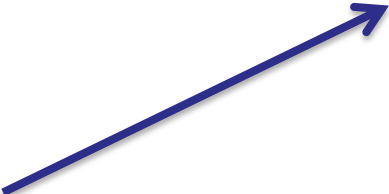
- df_t is the document frequency of t , namely the number of documents that contain term t
- Higher idf_t means that the term appears less often in the collection, and is therefore more informative (important)
 - e.g., **JKU** has high idf, while **the** has very low idf

Term weightings


- $\text{tf-idf}_{t,d}$ term weighting is the product of $\text{tf}_{t,d}$ and idf_t

$$x_{t,d} = \text{tf-idf}_{t,d} = \log(1 + \text{tc}_{t,d}) \times \log(M / \text{df}_t)$$

increases with the number of
occurrences within a document



increases with the rarity of the
term in the collection



- A well-known term weighting method!

Wrap-up!

- Use any term weightings to create document-term matrix

	$t1$	$t2$...	tN	sentiment
$d1$	$x_{t1,d1}$	$x_{t2,d1}$...	$x_{tN,d1}$	y_{d1}
$d2$	$x_{t1,d2}$	$x_{t2,d2}$...	$x_{tN,d2}$	y_{d2}
...
dM	$x_{t1,dM}$	$x_{t2,dM}$...	$x_{tN,dM}$	y_{dM}

- The rest is standard machine learning!
 - Model training
 - Hyper-parameter tuning and model selection
 - Evaluation

Agenda

- Introduction to Machine Learning
- Sentiment Analysis
- Feature Extraction
- ***Breaking* the curse of dimensionality!**

Supervised Sentiment Analysis

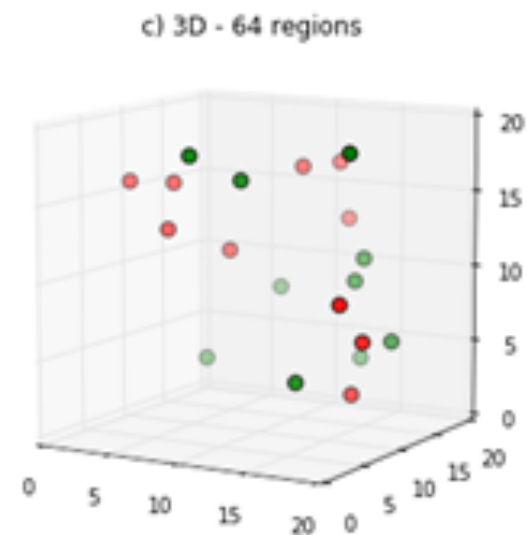
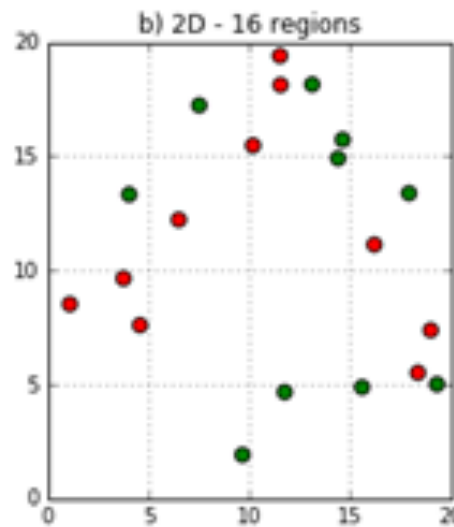
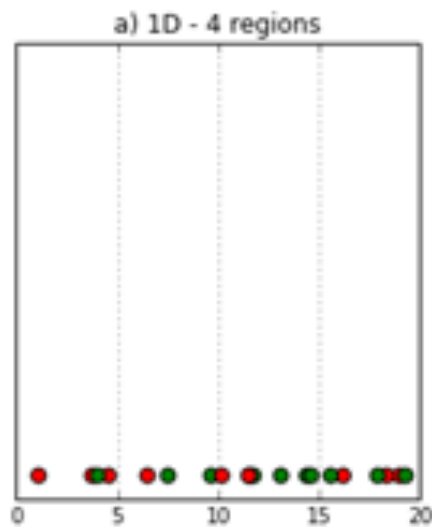
	$t1$	$t2$...	tN	sentiment
$d1$	$x_{t1,d1}$	$x_{t2,d1}$...	$x_{tN,d1}$	y_{d1}
$d2$	$x_{t1,d2}$	$x_{t2,d2}$...	$x_{tN,d2}$	y_{d2}
...
dM	$x_{t1,dM}$	$x_{t2,dM}$...	$x_{tN,dM}$	y_{dM}

$N \sim [20K - 500K]$
 $M \sim [10K - 100K]$

- The feature vectors are
 - **sparse** (a lot zeros)
 - in a very high dimension

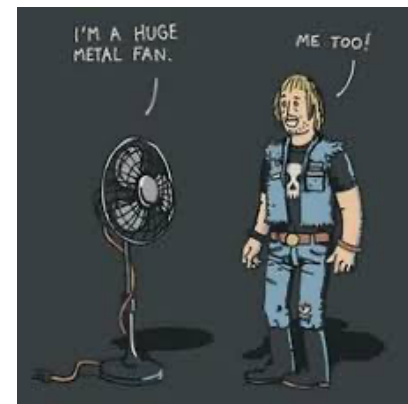
Curse of dimensionality

- Curse of dimensionality happens when the amount of data does not suffice to support the sparsity in dimensionality
- It causes
 - Data sparsity
 - Issues in measuring “closeness”



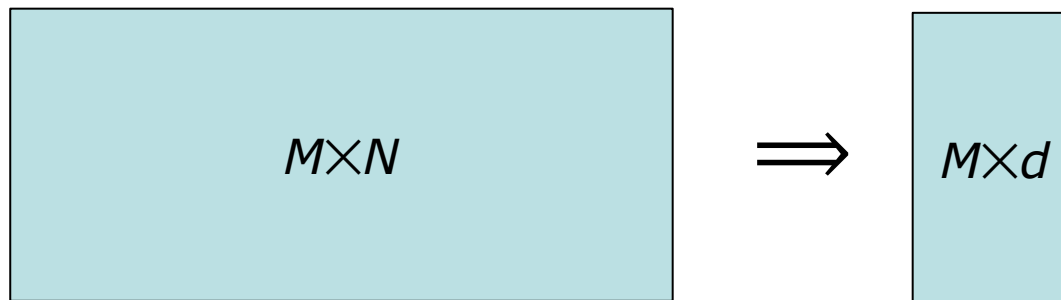
Curse of dimensionality

- Why low-dimensional vectors?
 - Easier to store and load (efficiency)
 - More efficient when used as features in ML models
 - Better generalization since the noise in data is reduced
 - Able to capture higher-order relations:
 - Synonymy like *car* and *automobile* can be merged into same dimensions
 - Polysomy like *bank (financial institution)* and *bank (bank of river)* can be separated into different dimensions



Feature (Dimensionality) reduction

- Feature selection
 - keep some important features and get rid of the rest!
- Dimensionality reduction
 - project data from high to a low dimensional space



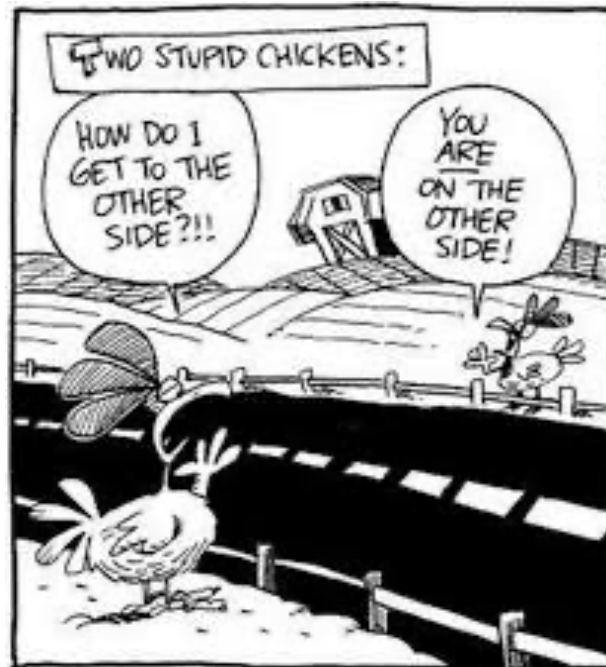
Feature selection

- During pre-processing
 - Remove **stop words** or very common words
 - tf-idf do it in a “soft” way *why?*
 - Remove very **rare words**
 - Usually done when creating dictionary
 - Stemming & lemmatization
- Features definition
 - Use only the words in a domain-specific **lexicon** as features
- Post-processing
 - Keep important features using some **informativeness** measures
 - Subset selection

Dimensionality reduction with LSA

- Latent Semantic Analysis (LSA)
 - A common method in Information Retrieval to capture semantics
 - Based on Singular Value Decomposition (SVD)

Semantics matters!



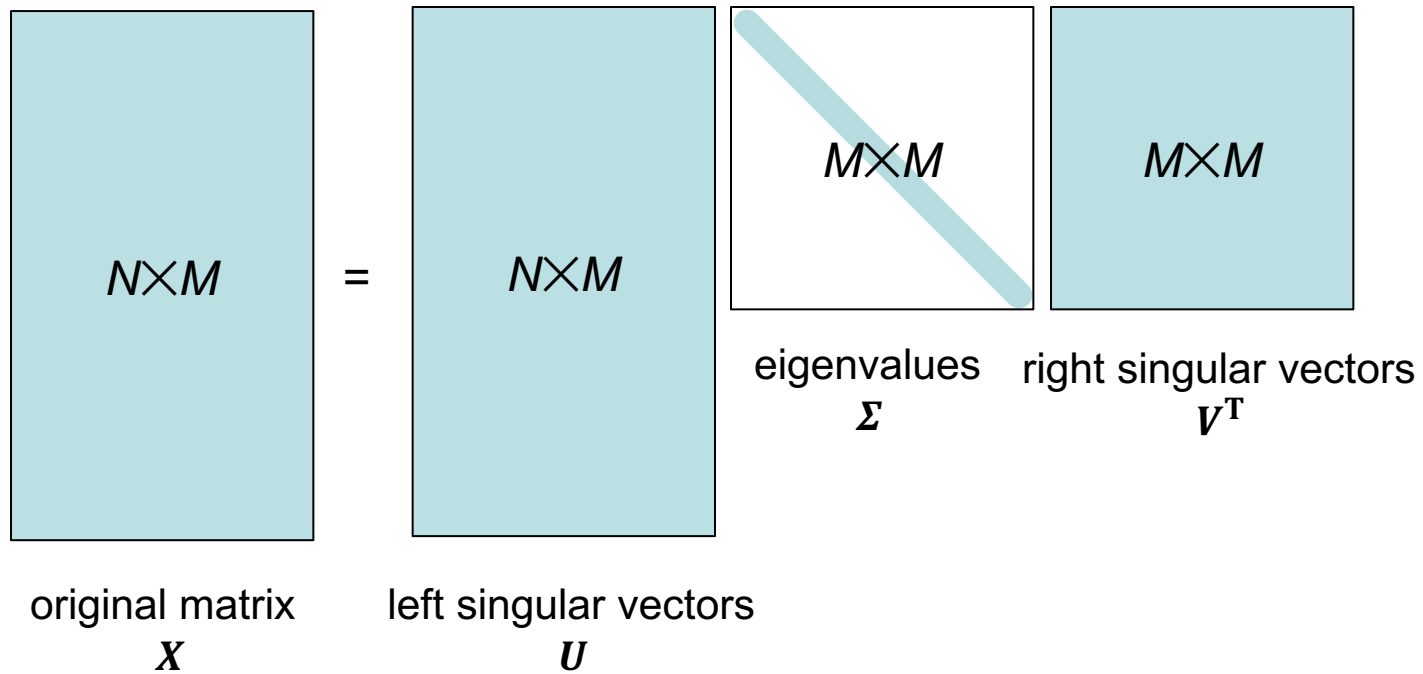
Singular Value Decomposition

- An $N \times M$ matrix X can be factorized to three matrices:

$$X = U\Sigma V^T$$

- U (left singular vectors) is an $N \times M$ unitary matrix
- Σ is an $M \times M$ diagonal matrix, diagonal entries
 - are eigenvalues,
 - show the importance of corresponding M dimensions in X
 - are all positive and sorted from large to small values
- V^T (right singular vectors) is an $M \times M$ unitary matrix

Singular Value Decomposition

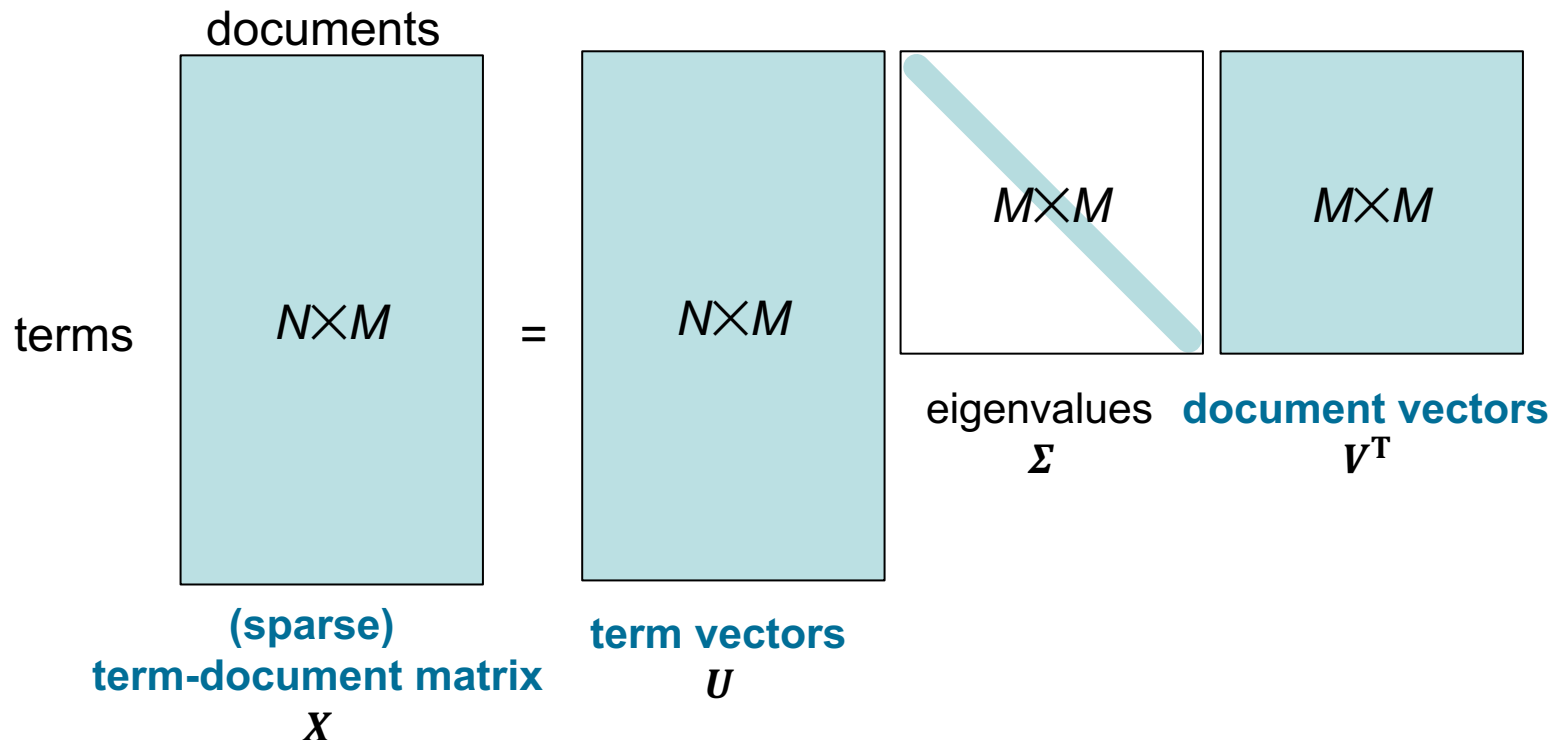


Latent Semantic Indexing

At Training Time

- Step 1: apply SVD on the **term-document** matrix of training data

☞ Not the document-term matrix! Although it is also possible to start with the document-term matrix, we follow the typical definition!



Latent Semantic Indexing

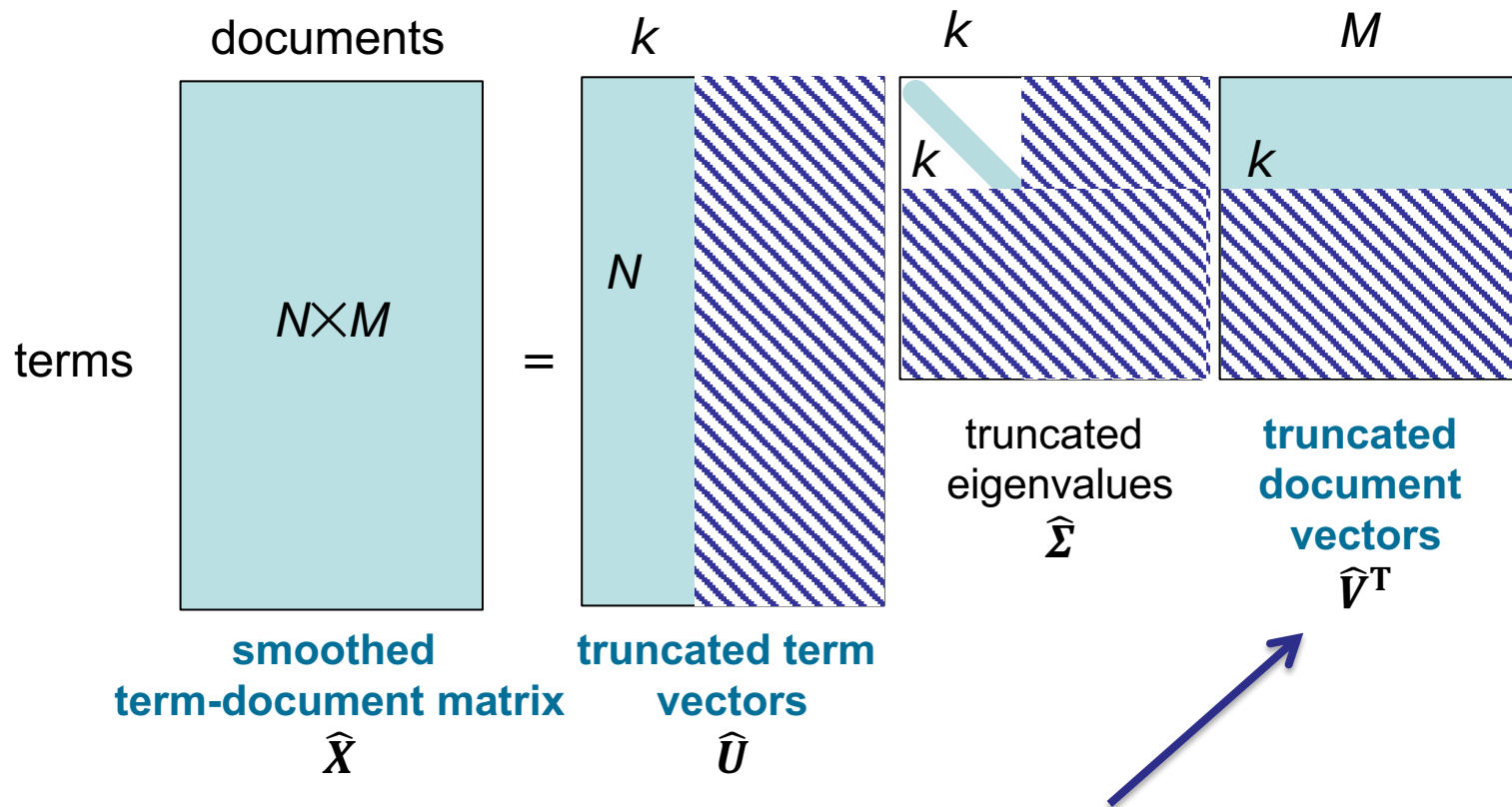
At Training Time

- Step 2: keep only top k eigenvalues in Σ and set the rest to zero, called $\hat{\Sigma}$
- Truncate the U and V^T matrices based on the changes in Σ , called \hat{U} and \hat{V}^T respectively
- If we multiply the truncated matrices, we have the rank k **least-squares approximation** to the original matrix

$$\hat{X} = \hat{U}\hat{\Sigma}\hat{V}^T$$

Latent Semantic Indexing

At Training Time



- \hat{V} matrix is the dense low-dimensional document vectors



Latent Semantic Indexing

At Inference Time (Validation or Test)

- Given a high-dimensional document vector \mathbf{d} in $N \times 1$ dimensions, we want to project it to the low-dimensional space, resulting in a new vector $\hat{\mathbf{d}}$ with $k \times 1$ dimensions
- done through this calculation:

$$\hat{\mathbf{d}} = \hat{\Sigma}^{-1} \hat{U}^T \mathbf{d}$$

- Exercise: calculate if the chain of dot products from \mathbf{d} to $\hat{\mathbf{d}}$ results to the correct dimension