

CS 2750 Machine Learning

Lecture 3

Designing a learning system II

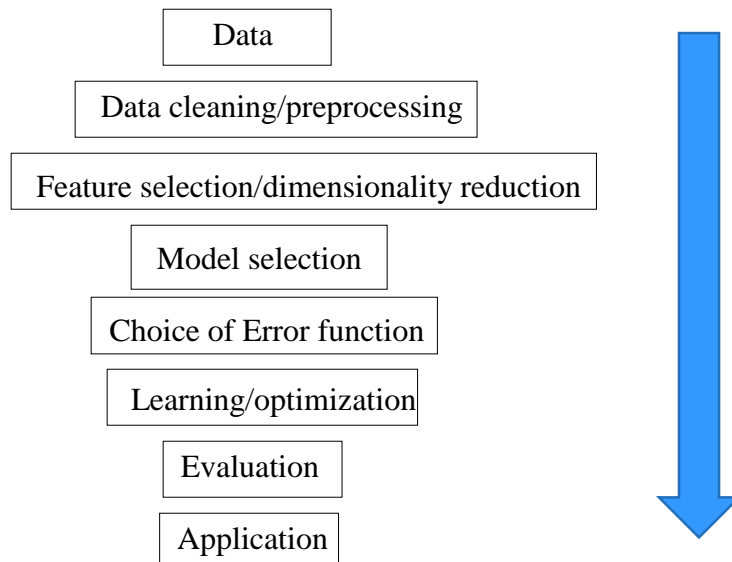
Milos Hauskrecht

milos@cs.pitt.edu

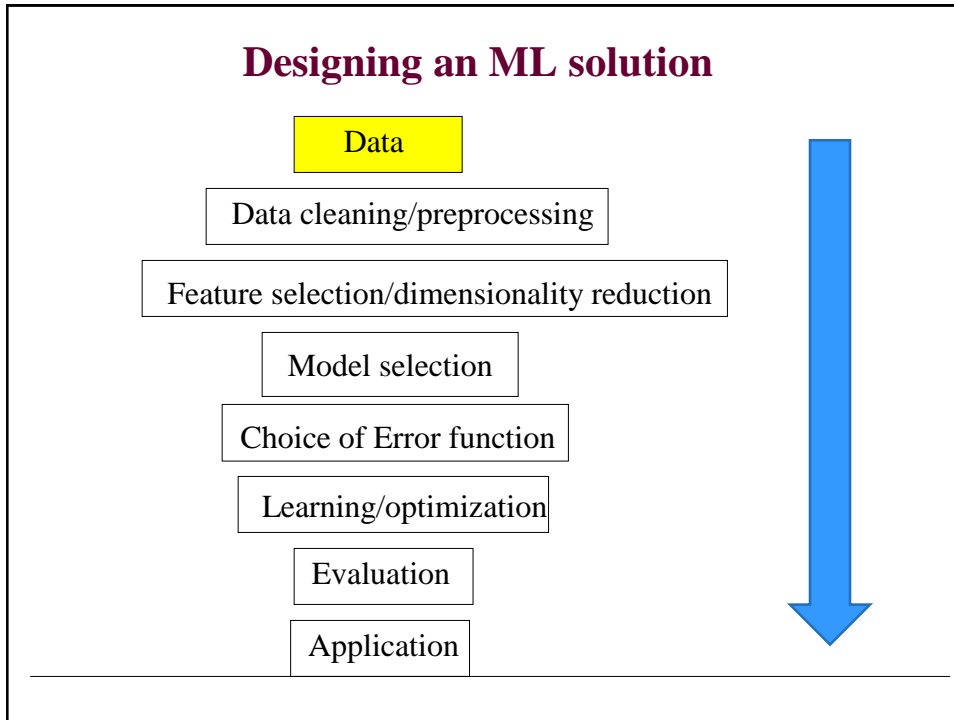
5329 Sennott Square, x4-8845

people.cs.pitt.edu/~milos/courses/cs2750-Spring2020/

Designing an ML solution



Designing an ML solution



Data source and data biases

- **Understand the data source**
- **Understand the data your models will be applied to**
- **Watch out for data biases:**
 - Make sure the data we make conclusions on are the same as data we used in the analysis
 - It is very easy to derive “unexpected” results when data used for analysis and learning are biased
- **Results (conclusions) derived for a biased dataset do not hold in general !!!**

Data biases

Example: Assume you want to build an ML program for predicting the stock behavior and for choosing your investment strategy

Data extraction:

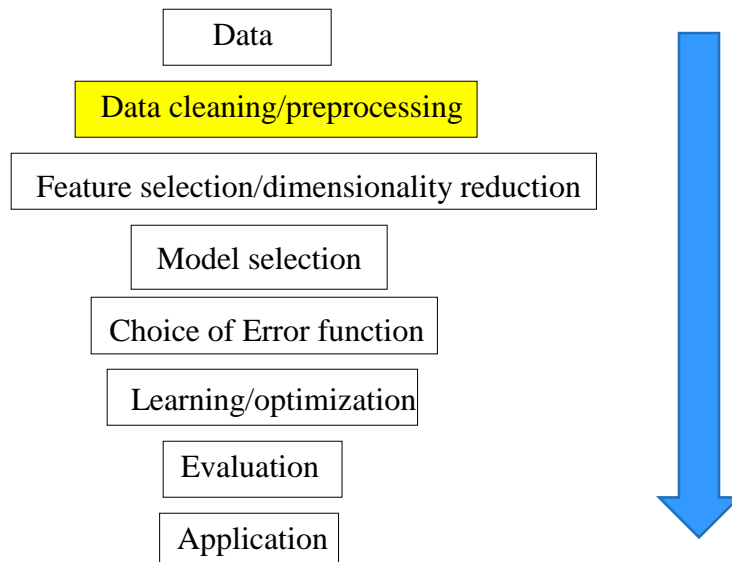
- pick companies that are traded on the stock market on January 2017
- Go back 30 years and extract all the data for these companies
- Use the data to build an ML model supporting your future investments

Question:

- **Would you trust the model?**
- **Are there any biases in the data?**

CS 2750 Machine Learning

Steps taken when designing an ML system



Data cleaning and preprocessing

Data you receive may not be perfect:

- Cleaning
- Preprocessing (conversions)

Cleaning:

- Get rid of errors, noise, outliers
- Removal of redundancies

Preprocessing:

- Renaming
- Rescaling (normalization)
- Discretizations
- Abstraction
- Aggregation
- New attributes

CS 2750 Machine Learning

Data preprocessing

Renaming (relabeling) categorical values to numbers

- dangerous in conjunction with some learning methods
- numbers will impose an order that is not warranted

Example:

- assume the following encoding of values High, Normal, Low

High \rightarrow 2

Normal \rightarrow 1

Low \rightarrow 0

- **2 > 1 implies High > Normal: Is it OK?**
- **1 > 0 implies Normal > Low: Is it OK?**
- **2 > 0 implies High > Low: Is it OK?**

?

Data preprocessing

Renaming (relabeling) categorical values to numbers

- dangerous in conjunction with some learning methods
- numbers will impose an order that is not warranted

Example:

- assume the following encoding of values High, Normal, Low

High \rightarrow 2

Normal \rightarrow 1

Low \rightarrow 0

- **2 > 1 implies High > Normal: Is it OK?**
- **1 > 0 implies Normal > Low: Is it OK?**
- **2 > 0 implies High > Low: Is it OK?**



Data preprocessing

Renaming (relabeling) categorical values to numbers

- dangerous in conjunction with some learning methods
- numbers will impose an order that is not warranted

High \rightarrow 2

Normal \rightarrow 1

Low \rightarrow 0

True \rightarrow 2

False \rightarrow 1

Unknown \rightarrow 0



Data preprocessing

Renaming (relabeling) categorical values to numbers

- dangerous in conjunction with some learning methods
- numbers will impose an order that is not warranted

High → 2
Normal → 1 ✓
Low → 0

True → 2
False → 1 ✗
Unknown → 0

Red → 2
Blue → 1 ?
Green → 0

Data preprocessing

Renaming (relabeling) categorical values to numbers

- dangerous in conjunction with some learning methods
- numbers will impose an order that is not warranted

High → 2
Normal → 1 ✓
Low → 0

True → 2
False → 1 ✗
Unknown → 0

Red → 2 ✗
Blue → 1 ✗
Green → 0

Data preprocessing

Renaming (relabeling) categorical values to numbers

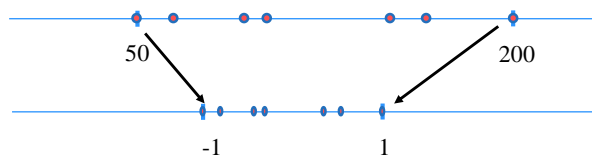
Problem: How to safely represent the different categories as numbers when no order exists?

Solution:

- Use indicator vector (or one-hot) representation.
- **Example: Red, Blue, Green colors**
 - 3 categories \rightarrow use a vector of size 3 with binary values
 - Encoding:
 - **Red:** (1,0,0);
 - **Blue:** (0,1,0);
 - **Green:** (0,0,1)

Data preprocessing

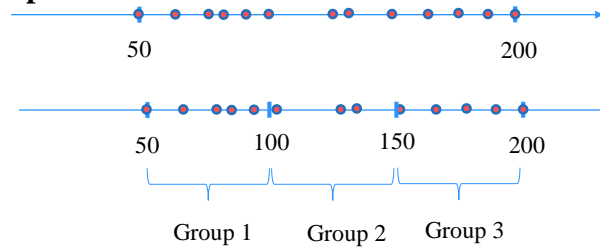
- **Rescaling (normalization):** continuous values transformed to some range, typically $[-1, 1]$ or $[0,1]$.



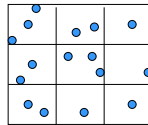
- Why normalization?
 - Some learning algorithms are sensitive to the values recorded in the specific input field and its magnitude

Data preprocessing

- **Discretization (binning):** continuous values to a finite set of discrete values
- **Example:**



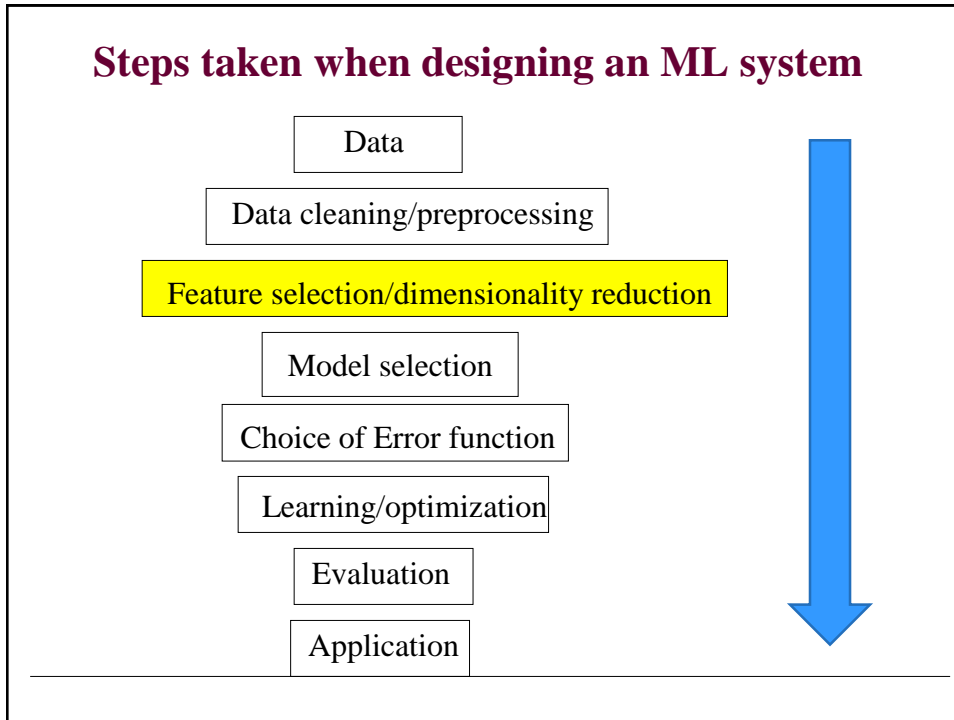
- **Example 2:**



Data preprocessing

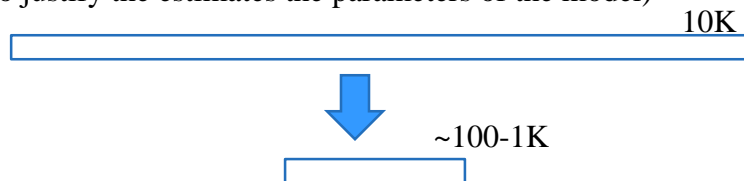
- **Abstraction:** merge together categorical values
- **Aggregation:** summary or aggregation operations, such minimum value, maximum value, average etc.
- **New attributes:**
 - example: obesity-factor = weight/height

Steps taken when designing an ML system



Feature selection/dimensionality reduction

- **The size (dimensionality) of an instance** can be enormous
 $x_i = (x_i^1, x_i^2, \dots, x_i^d)$ d - very large
- **Problem:** Too many parameters to learn (not enough samples to justify the estimates the parameters of the model)



Example: **document classification**

- 10,000 different words
- Big vector: counts of occurrences of different words

Feature selection/dimensionality reduction

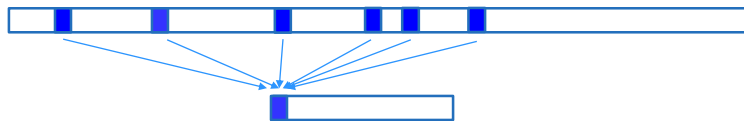
- Dimensionality reduction solutions

- Extract a small subset of original inputs

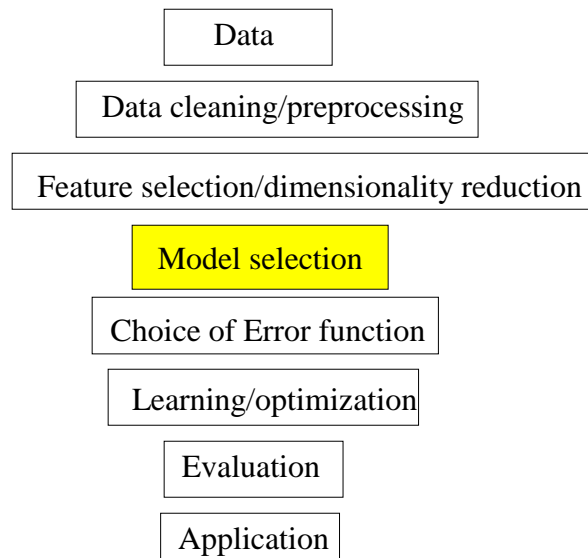


- Project inputs into a lower dimensional vector:

- PCA – principal component analysis
- Auto-encoders
- Latent variable models



Steps taken when designing an ML system

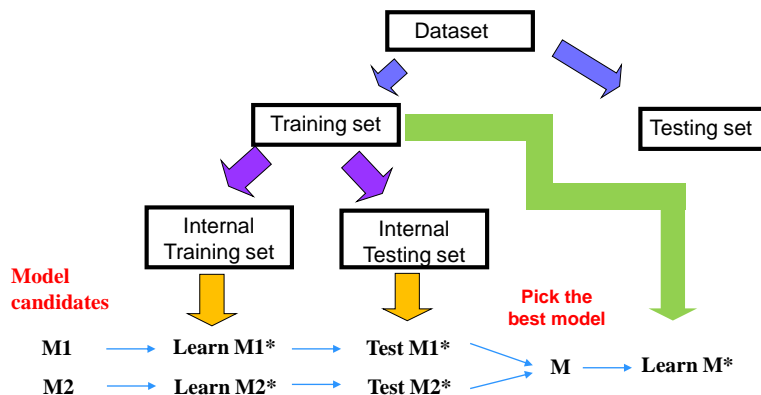


Model selection

- **What is the right model to learn?**
 - A prior knowledge helps a lot, but still a lot of guessing
 - Initial data analysis and visualization
 - We can make a good guess about the form of the distribution, shape of the function by looking at data
 - Independences and correlations
- **Overfitting problem**
 - Take into account the **bias and variance** of error estimates
 - Simpler (more biased) model – parameters can be estimated more reliably (smaller variance of estimates)
 - Complex model with many parameters – parameter estimates are less reliable (large variance of the estimate)

Solutions for overfitting

- A. Use internal train and test splitting.** Basically, hold some data out of the training set (called validation set) to decide on the model first, then train the picked model

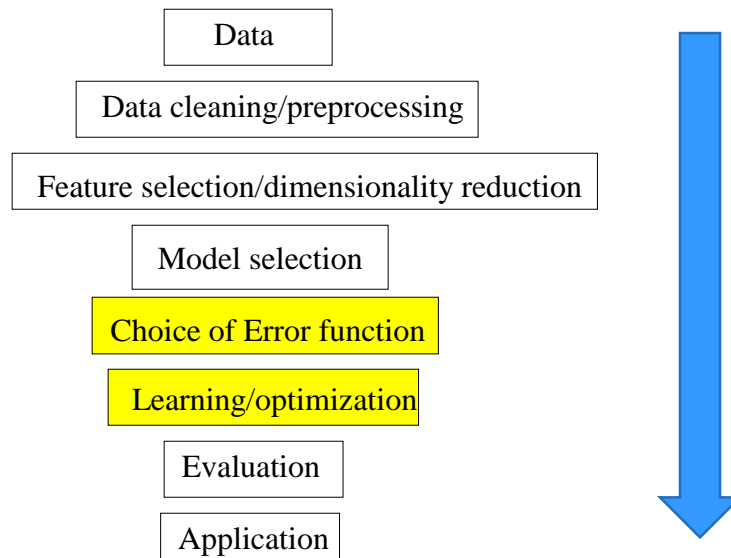


Solutions for overfitting

How to decide among models with different complexity?

- **Regularization (Occam's Razor)**
 - Penalize for the model complexity (number of parameters) in the objective function
 - Lasso or Ridge regularization
 - Explicit preference towards simple models

Steps taken when designing an ML system



Learning: objective functions

- **Learning = optimization problem.** Various criteria:

- **Mean square error**

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} Error(\mathbf{w}) \quad Error(\mathbf{w}) = \frac{1}{N} \sum_{i=1, \dots, N} (y_i - f(x_i, \mathbf{w}))^2$$

- **Maximum likelihood (ML) criterion**

$$\Theta^* = \max_{\Theta} P(D | \Theta) \quad Error(\Theta) = -\log P(D | \Theta)$$

- **Maximum posterior probability (MAP)**

$$\Theta^* = \max_{\Theta} P(\Theta | D) \quad P(\Theta | D) = \frac{P(D | \Theta)P(\Theta)}{P(D)}$$

CS 2750 Machine Learning

Learning

Learning = optimization problem

- Optimization problems can be hard to solve. Right choice of a model and an error function makes a difference.
- **Parameter optimizations**

- Gradient descent, Conjugate gradient
- Newton-Raphson
- Levenberg-Marquard

Some can be carried **on-line** on a sample by sample basis

Combinatorial optimizations (over discrete spaces):

- Hill-climbing
- Simulated-annealing
- Genetic algorithms

CS 2750 Machine Learning

Parametric optimizations

- Sometimes can be solved directly but this depends on the objective function and the model
 - **Example:** squared error criterion for the linear regression
- Very often the objective function to be optimized is not that nice.

$$Error(\mathbf{w}) = f(\mathbf{w}) \quad \mathbf{w} = (w_0, w_1, w_2 \dots w_k)$$

- a complex function of weights (parameters)

$$\text{Goal: } \mathbf{w}^* = \arg \min_{\mathbf{w}} f(\mathbf{w})$$

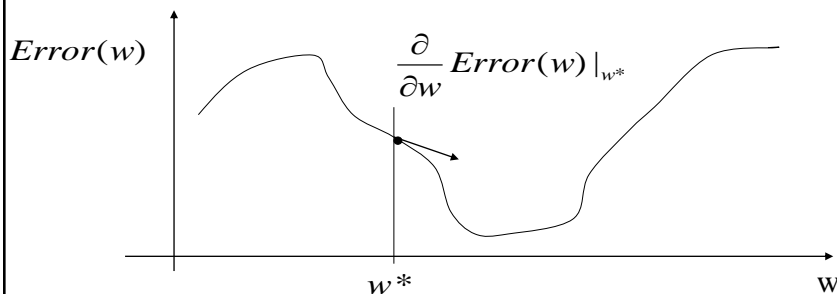
- One solution: **iterative optimization methods**
- **Example: Gradient-descent method**

Idea: move the weights (free parameters) gradually in the error decreasing direction

CS 2750 Machine Learning

Gradient descent method

- Descend to the minimum of the function using the gradient information

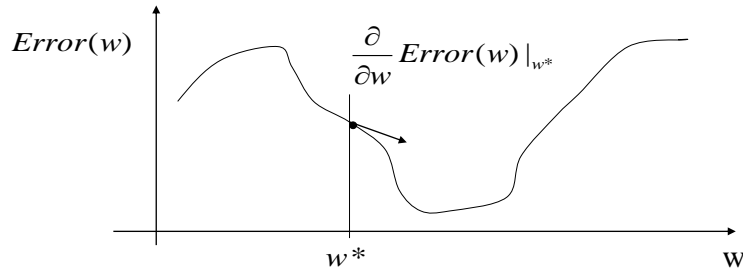


- Change the parameter value of w according to the gradient

$$w \leftarrow w^* - \alpha \frac{\partial}{\partial w} Error(w) |_{w^*}$$

CS 2750 Machine Learning

Gradient descent method



- New value of the parameter

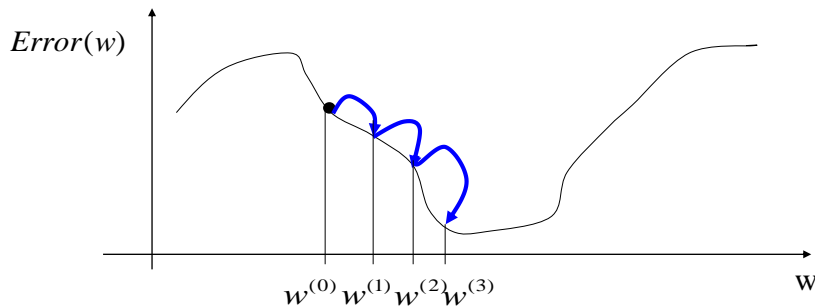
$$w \leftarrow w^* - \alpha \frac{\partial}{\partial w} Error(w) |_{w^*}$$

$\alpha > 0$ - a learning rate (scales the gradient changes)

CS 2750 Machine Learning

Gradient descent method

- To get to the function minimum repeat (iterate) the gradient based update few times



- **Problems:** local optima, saddle points, slow convergence
- More complex optimization techniques use additional information (e.g. second derivatives)

CS 2750 Machine Learning

On-line learning (optimization)

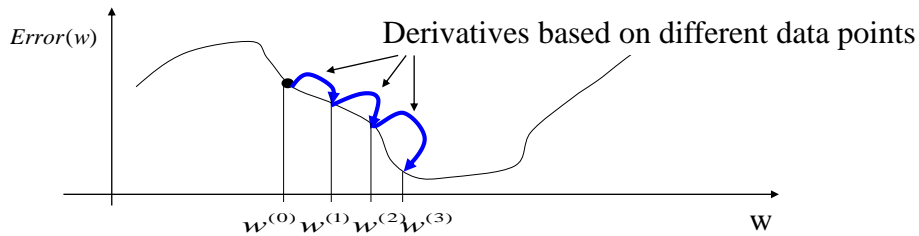
- Error function looks at all data points at the same time

$$\text{E.g. } Error(\mathbf{w}) = \frac{1}{n} \sum_{i=1, \dots, n} (y_i - f(x_i, \mathbf{w}))^2$$

- **On-line error** - separates the contribution from a data point

$$Error_{\text{ON-LINE}}(\mathbf{w}) = (y_i - f(x_i, \mathbf{w}))^2$$

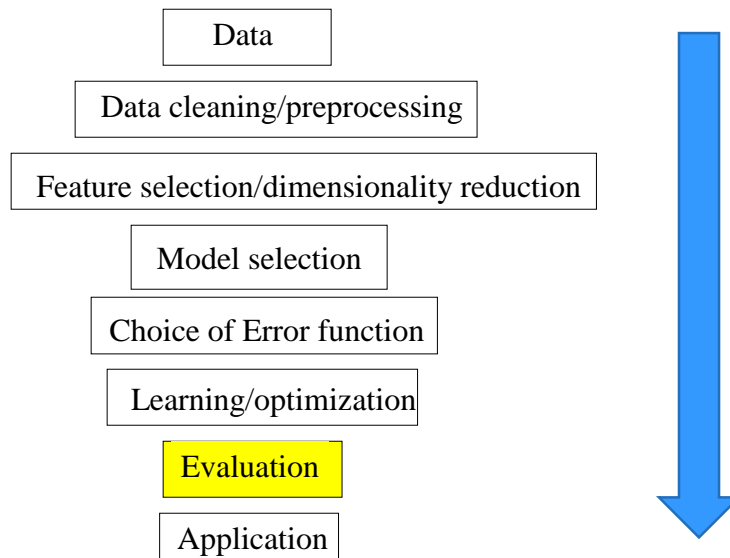
- **Example: On-line gradient descent**



- **Advantages:** 1. simple learning algorithm
2. no need to store data (on-line data streams)

CS 2750 Machine Learning

Steps taken when designing an ML system



Evaluation measures

Regression model $f: X \rightarrow Y$ where Y is real valued

- The error is calculated on the data D , say

$$MSE(D, f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

- This is an estimate of the error for f on the complete population

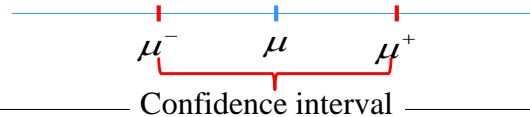
Important question:

- How close is our estimate to the true mean error?

To answer the question we need to resort to statistics:

- How confident we are the true error falls into interval around our estimate μ ?

Answer: with probability 0.95 the true error is in interval $[\mu^-, \mu^+]$

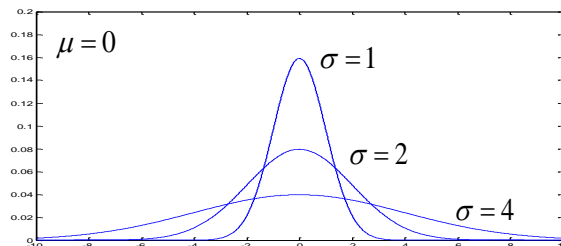


Evaluation

- Central limit theorem:**

Let random variables X_1, X_2, \dots, X_n form a random sample from a **distribution** with mean μ and variance σ , then if the sample n is large, the distribution

$$\sum_{i=1}^n X_i \approx N(n\mu, n\sigma^2) \quad \text{or} \quad \frac{1}{n} \sum_{i=1}^n X_i \approx N(\mu, \sigma^2 / n)$$



Statistical significance test

- **Statistical tests for the mean**

- **H0 (null hypothesis)**

$$E[X] = \mu^0$$

- **H1 (alternative hypothesis)**

$$E[X] \neq \mu^0$$

- **Basic idea:**

we use the sample mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

and check how probable it is that $E[X] = \mu^0$ holds

If the probability that \bar{X} comes from the normal distribution with mean μ^0 is small – we reject the null hypothesis on that probability level

Statistical significance test

- **Statistical tests for the mean**

- **H0 (null hypothesis)**

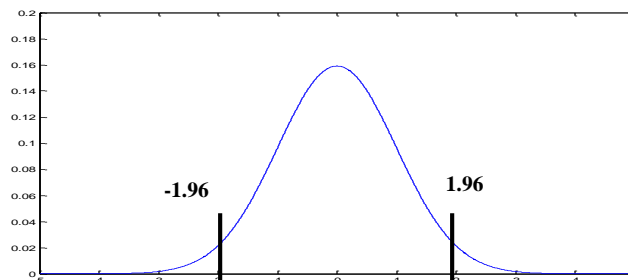
$$E[X] = \mu^0$$

- **H1 (alternative hypothesis)**

$$E[X] \neq \mu^0$$

- **Assume we know the standard deviation σ for the sample**

$$z = \frac{\bar{X} - \mu^0}{\sigma} \sqrt{n} \approx N(0,1) \quad \text{with} \quad P=0.95 \quad z \in [-1.96, 1.96]$$

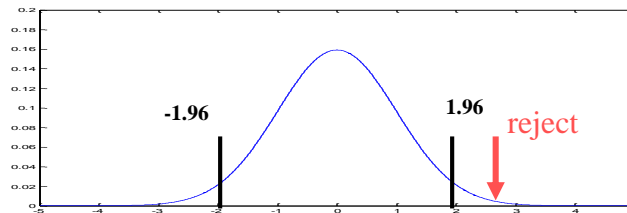


Statistical significance test

- **Statistical tests for the mean** $E[X] = \mu^0$
 - **H0 (null hypothesis)**
- **Assume we know the standard deviation σ**

$$z = \frac{\bar{X} - \mu^0}{\sigma} \sqrt{n} \approx N(0,1) \quad \text{with} \quad P=0.95 \quad z \in [-1.96, 1.96]$$

- **Z-test: If z is outside of the interval – reject the null hypothesis at significance level $(1 - P)$ if $P=0.95$ it is 0.05**



Statistical significance test

- **Statistical tests for the mean** $E[X] = \mu^0$
 - **H0 (null hypothesis)**
- **Problem: we do not know the standard deviation σ**

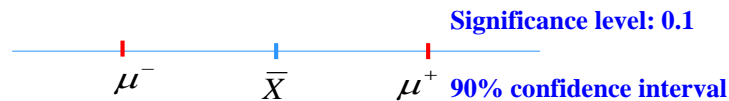
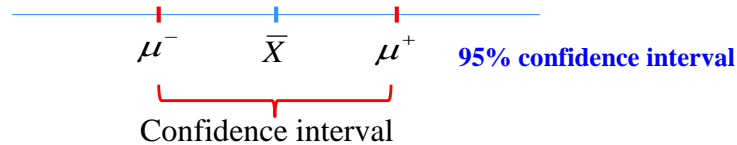
- **Solution:** $t = \frac{\bar{X} - \mu^0}{s} \sqrt{n} \approx t\text{-distribution (Student distribution)}$

$$s = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}} \quad \text{- Estimate of the standard deviation}$$

- **T-test: If t is outside of the tabulated interval reject the null hypothesis at the corresponding significance level**

Confidence interval

- Assume we have calculated the average error $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
- There are many values of μ^0 around it that are not rejected at some significance level (say 0.05)
- These values form a confidence interval around it



Statistical tests

The statistical tests lets us answer:

- The probability with which the true error falls into the interval around our estimate, say :

$$MSE(D, f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

- Compare two models M1 and M2 and determine based on the error on the data entries the probability with which model M1 is different (or better) than M2

$$MSE(D, f_1) = \frac{1}{n} \sum_{i=1}^n (y_i - f_1(x_i))^2 \quad MSE(D, f_2) = \frac{1}{n} \sum_{i=1}^n (y_i - f_2(x_i))^2$$

Trick:

$$\begin{aligned}
 MSE(D, f_1) - MSE(D, f_2) &= \frac{1}{n} \sum_{i=1}^n (y_i - f_1(x_i))^2 - \frac{1}{n} \sum_{i=1}^n (y_i - f_2(x_i))^2 \\
 &= \frac{1}{n} \sum_{i=1}^n (y_i - f_1(x_i))^2 - (y_i - f_2(x_i))^2
 \end{aligned}$$

Evaluation measures for classification

Assume binary classification:

- **Confusion matrix** represents all possible combination of true and predicted values

		Actual	
		Case	Control
Prediction	Case	TP 0.3	FP 0.1
	Control	FN 0.2	TN 0.4

TP – true positive
 FP – false positive
 TN – true negative
 FN – false negative

Evaluation measures for classification

Evaluation stats calculated from the confusion matrix:

		Actual	
		Case	Control
Prediction	Case	TP 0.3	FP 0.1
	Control	FN 0.2	TN 0.4

TP – true positive
 FP – false positive
 TN – true negative
 FN – false negative

Misclassification error:

$$E = FP + FN$$

Accuracy:

$$Accuracy = TP + TN$$

Sensitivity:

$$SN = \frac{TP}{TP + FN}$$

Specificity:

$$SP = \frac{TN}{TN + FP}$$

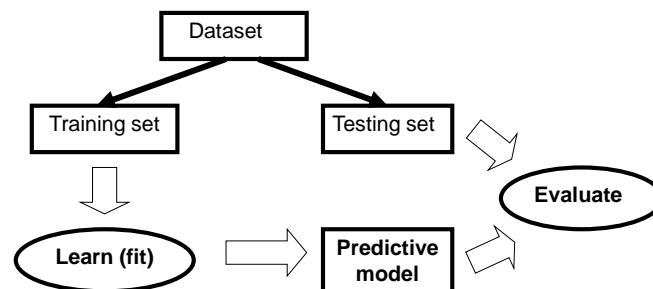
Evaluation measures for classification

More complex classification evaluation measures:

- AUROC:
 - Area under the Receiver operating curve
- AUPRC
 - Area under the Precision-Recall curve

Evaluation of models

- We started with a simple holdout method



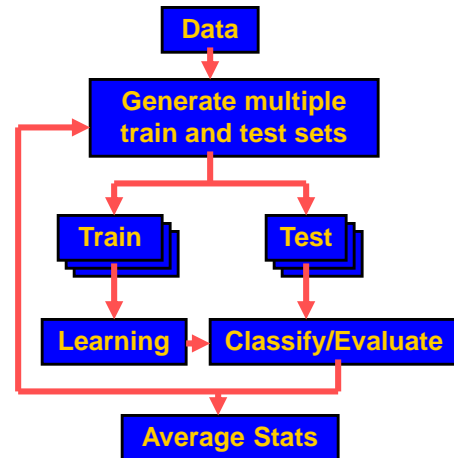
Problem: the mean error results may be influenced by a lucky or an unlucky **training and testing** split especially for a small size D

Solution: try multiple train-test splits and average their results

Evaluation of models via random resampling

Other more complex methods

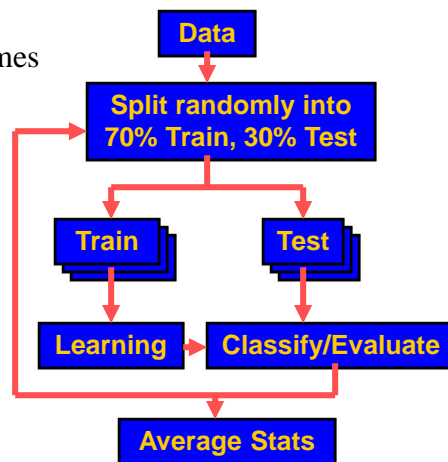
- Use multiple train/test sets
- Based on various random re-sampling schemes:
 - Random sub-sampling
 - Cross-validation
 - Bootstrap



CS 2750 Machine Learning

Evaluation of models using random subsampling

- Random sub-sampling
 - Repeat a simple holdout method k times

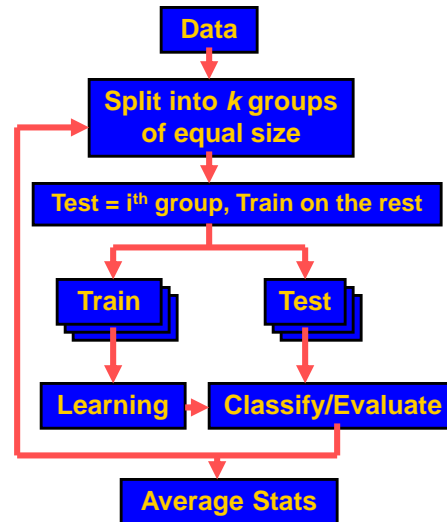


CS 2750 Machine Learning

Evaluation of models using k-fold cross-validation

Cross-validation (k-fold)

- Divide data into k disjoint groups, test on k-th group/train on the rest
- Typically 10-fold cross-validation
- Leave one out cross-validation (k = size of the data D)



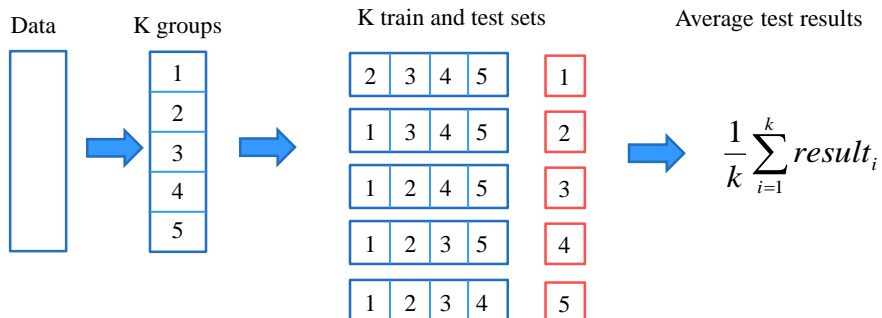
CS 2750 Machine Learning

Evaluation of models using k-fold cross-validation

Cross-validation (k-fold)

- Divide data into k disjoint groups,
- For every group i, test on i-th group and train on the rest
- Gives k models and k test results

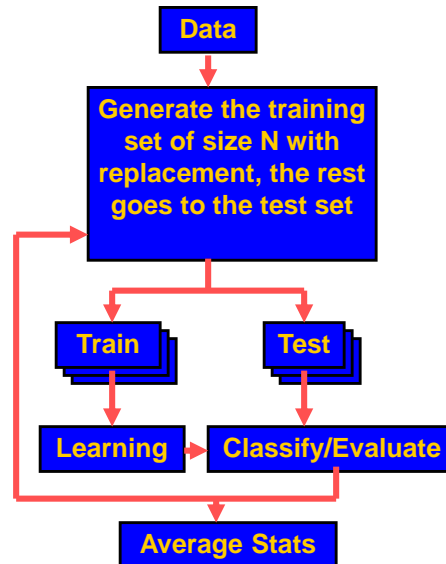
Example: k=5 (5-fold crossvalidation)



Evaluation of models using bootstrap

Bootstrap

- The training set of size N = size of the data D
- Sampling with the replacement



CS 2750 Machine Learning

Evaluation.

- **Simple holdout method.**
 - Divide the data to the training and test data.
- **Other more complex methods**
 - Based on cross-validation, random sub-sampling.
- What if we want to compare the predictive performance on a classification or a regression problem for two different learning methods?
- **Solution:** compare the error results on the test data set
- **Possible answer:** the method with better (smaller) testing error gives a better generalization error.
- But we need to use statistics to validate the choice

CS 2750 Machine Learning