



UiO • Department of informatics
University of Oslo

Adaptive Methods for *Lecture* Data-based Decision Making *5*

IN-STK 5000 / 9000

Autumn 2022

Slides by Dr. Anne-Marie George, UiO

What we talk about today

**Supervised
Learning:
Classification
& Regression**

**Quality
Measures for
Classification**

**Quality
Measures for
Regression**

**Anonymity (&
it's Problems)**

**Differential
Privacy**



**Guest Lecture:
Chinmayi
Baramashetru,
14.20 Fortress**

Classification

Task: Find a function (*classifier*) to classify some input

Training data: Examples of the form (data point x , label/class y_x)

DISCRETE SET OF POSSIBLE LABELS

Example:

Task: Label pictures according to what they show!

Training Data: Picture + Label (Dog or Blueberry muffin)



→ Find a classifier that gives for any (new) picture a correct label.

Regression

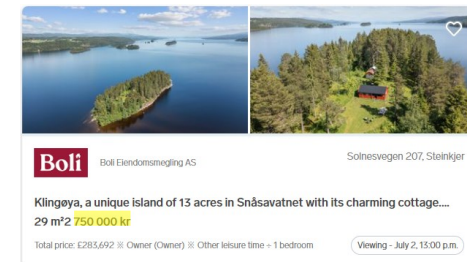
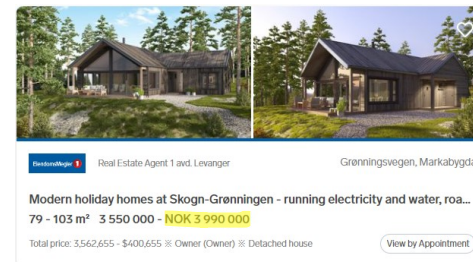
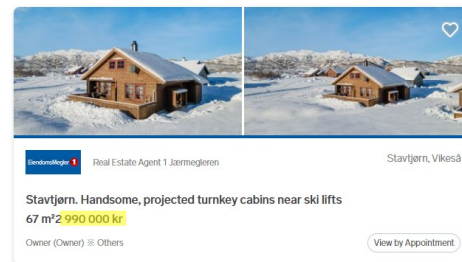
Task: Find a function (*regressor*) to evaluate some input

Training data: Examples of the form (data point x , label/value y_x)

Example: CONTINUOUS SET OF POSSIBLE LABELS

Task: Predicting housing prices!

Training Data: House (m^2 , post code, ...) + Price



→ Find a regressor that gives for any (new) house a correct price.

Classification & Regression

Decisions:

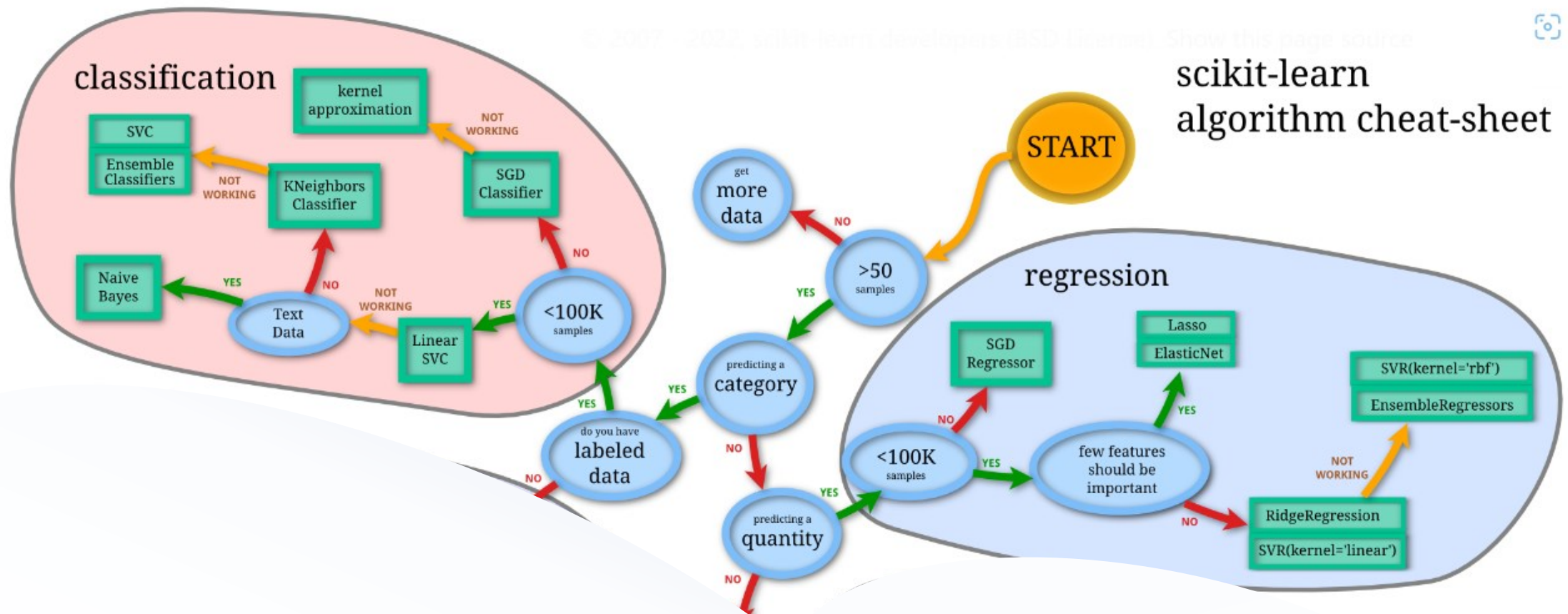
- What do we want to predict?
- What type of model do we need?
- What algorithm do we use?
- What kind of data do we need?
- How much data do we need?
- How do we measure the quality of the classifier/regressor?

Example: Estimate success chances of applicants to a study program!

- pass / fail, grades (A-F), ...
- binary classifier, regression, ...
- random forest, KNN, ...
- age, gender, nationality, GPA, ...
- Hundreds? Thousands? ...?
- accuracy, precision, ...?



scikit-learn algorithm cheat-sheet



Quality Measures for Classification

How do we measure quality of our predictions?



How to measure the quality of a classifier?



I build a classifier! 😊

→ It predicts the correct result in 91 % of the cases!

Is this a good classifier?

Depends:

- Is the dataset skewed?
- Are we doing better than, e.g., always predicting label1 no matter the input?
- What errors do we have?
- Where do they occur?

True / False Positives / Negatives

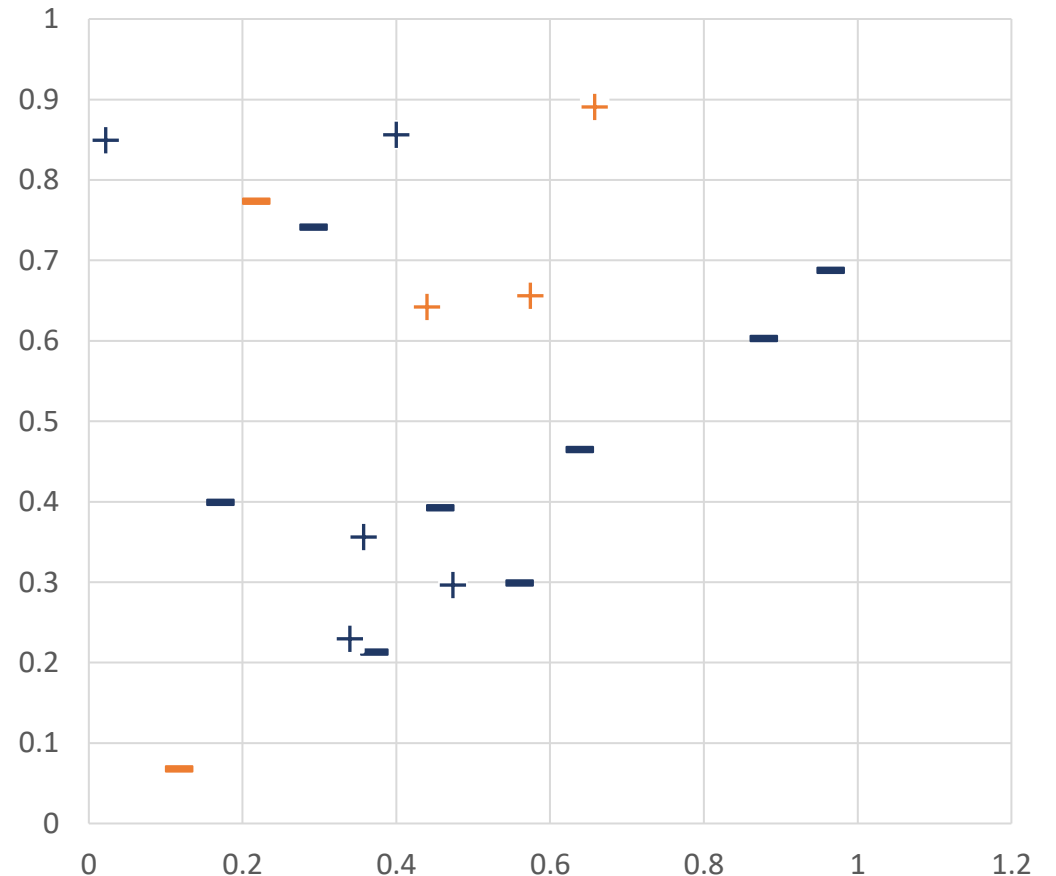
- Assumption:
There are only two labels +/- (or 1/0)
→ *Binary* classification

Pre- dicted Class	Actual Class		
		+	-
	+	True Positives (TP)	False Positives (FP)
	-	False Negatives (FN)	True Negatives (TN)

True / False

Positives / Negatives

- Assumption:
There are only two labels +/- (or 0/1)
→ *Binary* classification
- True Positives (TP) +
datapoints that are correctly classified as +
- False Positives (FP) +
datapoints that are falsely classified as +
- True Negatives (TN) -
datapoints that are correctly classified as -
- False Negatives (FN) -
datapoints that are falsely classified as -



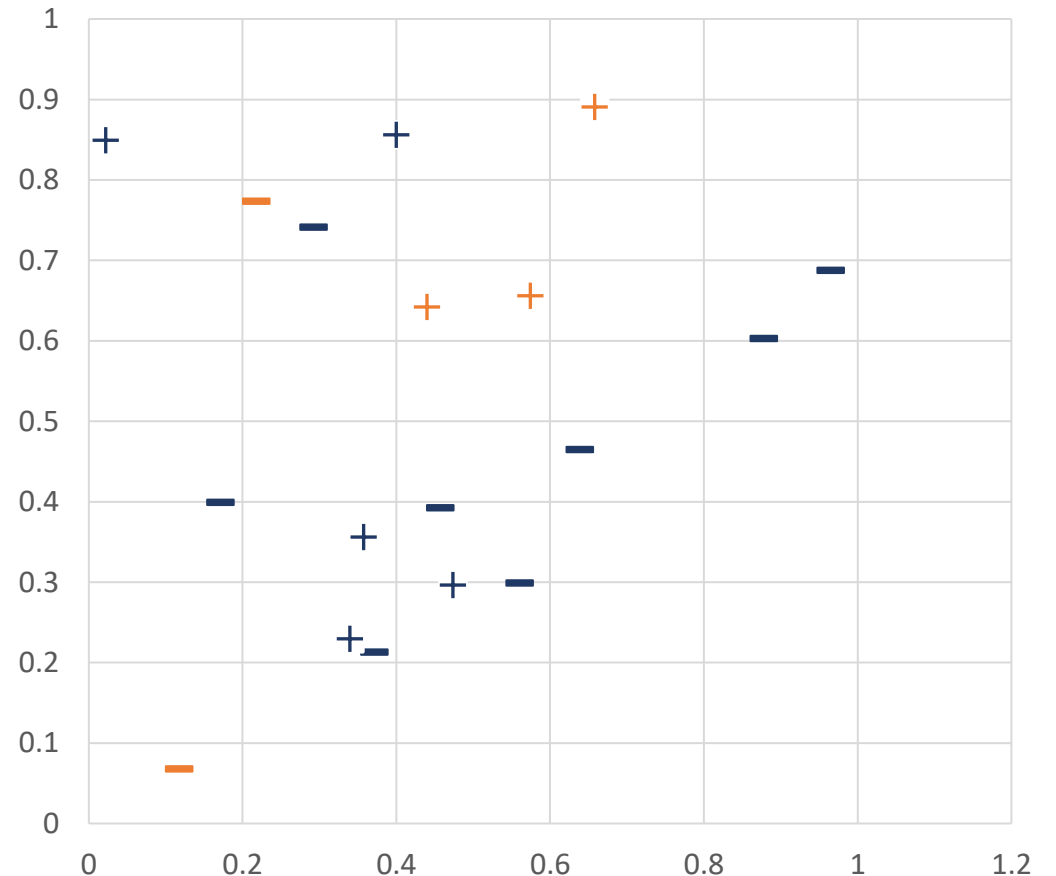
True / False

Positives / Negatives

- Example:
Predicting Covid:
positive or negative outcome.

What is the meaning of ... ?

- True Positives (TP) +
datapoints that are correctly classified as +
- False Positives (FP) +
datapoints that are falsely classified as +
- True Negatives (TN) -
datapoints that are correctly classified as -
- False Negatives (FN) -
datapoints that are falsely classified as -



Accuracy - for **binary** classification

- Accuracy: The fraction of correct predictions.

$$\begin{aligned} \text{Accuracy} &= \frac{\text{True Positives}(TP) + \text{True Negatives}(TN)}{\text{Total Number of Predictions}} \\ &= \frac{TP+TN}{TP+TN+FP+FN} \end{aligned}$$

Accuracy - for **multi-class** classification

- Accuracy: The fraction of correct predictions.

$$\textit{Accuracy} = \frac{\textit{Correct Predictions}}{\textit{Total Number of Predictions}}$$

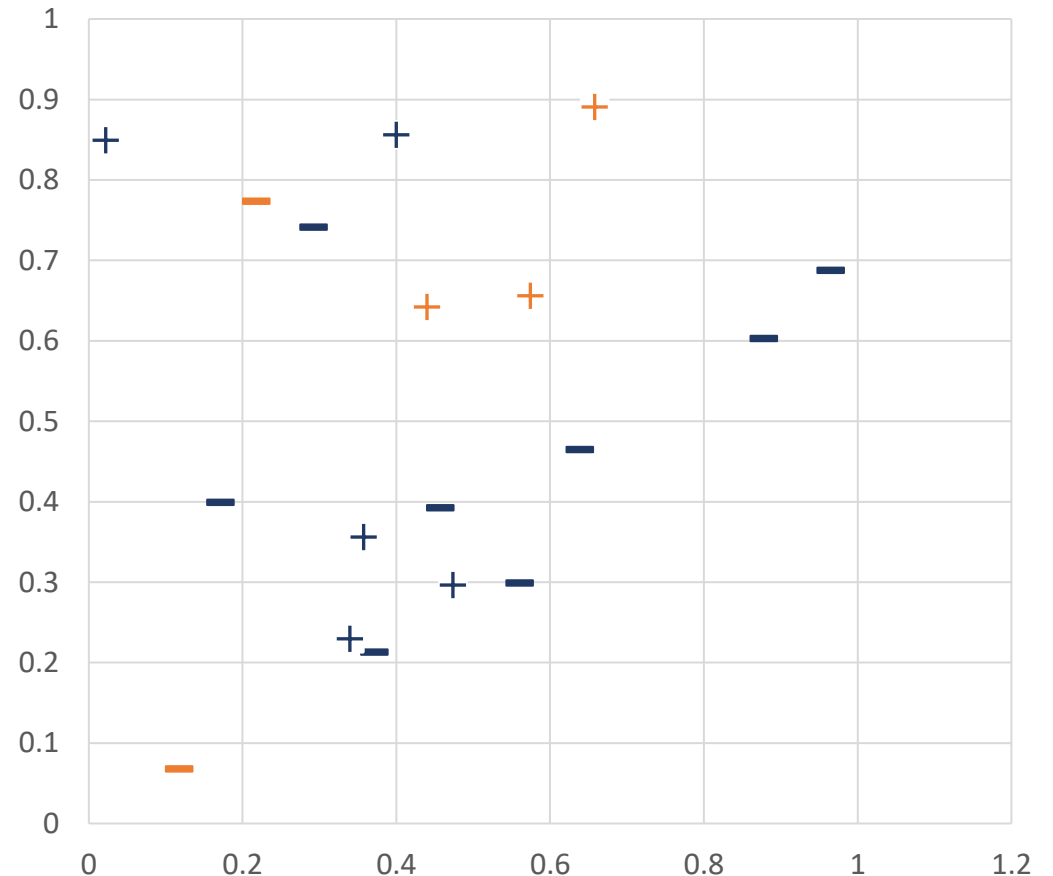
→ This *only* informs about the fraction of correct predictions, and *not* where the errors occur or what type(s) of error we have!

True / False

Positives / Negatives

- Example:
Predicting Covid:
positive or negative outcome.
→ What do we care about the most?

- True Positives (TP) +
datapoints that are correctly classified as +
- False Positives (FP) +
datapoints that are falsely classified as +
- True Negatives (TN) -
datapoints that are correctly classified as -
- False Negatives (FN) -
datapoints that are falsely classified as -



How to measure the quality of a classifier?



I improved my (Covid) classifier! 😊

- It has an accuracy of 0.95 (instead 0.91 from before)
- It has almost the same number of False Negatives (FN)!

Is this a better classifier?

- If the data is skewed, e.g., almost no one has Covid, then simply predicting “no Covid” more often increases the accuracy and doesn’t change the number of FN much (in expectation)!

How can we detect this?

Precision & Recall

- Precision: The fraction of true positive vs. all positive predictions.
(also: *positive predictive value*)

$$\text{Precision} = \frac{\text{True Positives}(TP)}{\text{Number of Positive Predictions}} = \frac{TP}{TP + FP}$$

- Recall: The fraction of true positive predictions vs. actual positives.
(also: *true positive rate* or *sensitivity*)

$$\text{Recall} = \frac{\text{True Positives}(TP)}{\text{Number of Actual Positives}} = \frac{TP}{TP + FN}$$

What is better?
High or low values?



Tradeoff: Precision & Recall

Example: Predicting risk to reoffend: “no crime” = +, “reoffend” = -

Given probability $h(x)$ of “reoffend” for any criminal x :

- If $h(x) < \theta = 0.8 \rightarrow$ Predict “no crime”, can be released.
- If $h(x) \geq \theta = 0.8 \rightarrow$ Predict “reoffend”, must stay in prison.



Safety of public: Predict “no crime” only with high confidence!

⚡ \rightarrow FP low, FN high: $\nearrow Precision = \frac{TP}{TP+FP}$ $\searrow Recall = \frac{TP}{TP+FN}$

Individual justice: Predict “reoffend” only with high confidence!

\rightarrow FP high, FN low: $\searrow Precision = \frac{TP}{TP+FP}$ $\nearrow Recall = \frac{TP}{TP+FN}$

Other Measures

→ A good measure is application dependent!

Sources: [5][6][7][8][9][10][11][12] view · talk · edit

		Predicted condition			
		Positive (PP)	Negative (PN)	Informedness, bookmaker informedness (BM) $= \text{TPR} + \text{TNR} - 1$	Prevalence threshold (PT) $= \frac{\sqrt{\text{TPR} \times \text{FPR}} - \text{FPR}}{\text{TPR} - \text{FPR}}$
Actual condition	Total population $= P + N$				
	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{\text{TP}}{P} = 1 - \text{FNR}$	False negative rate (FNR), miss rate $= \frac{\text{FN}}{P} = 1 - \text{TPR}$
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out $= \frac{\text{FP}}{N} = 1 - \text{TNR}$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{\text{TN}}{N} = 1 - \text{FPR}$
	Prevalence $= \frac{P}{P+N}$	Positive predictive value (PPV), precision $= \frac{\text{TP}}{\text{PP}} = 1 - \text{FDR}$	False omission rate (FOR) $= \frac{\text{FN}}{\text{PN}} = 1 - \text{NPV}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$
	Accuracy (ACC) $= \frac{\text{TP} + \text{TN}}{P + N}$	False discovery rate (FDR) $= \frac{\text{FP}}{\text{PP}} = 1 - \text{PPV}$	Negative predictive value (NPV) $= \frac{\text{TN}}{\text{PN}}$ $= 1 - \text{FOR}$	Markedness (MK), deltaP (Δp) $= \text{PPV} + \text{NPV} - 1$	Diagnostic odds ratio (DOR) $= \frac{\text{LR}+}{\text{LR}-}$
	Balanced accuracy (BA) $= \frac{\text{TPR} + \text{TNR}}{2}$	F ₁ score $= \frac{2\text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$	Fowlkes–Mallows index (FM) $= \sqrt{\text{PPV} \times \text{TPR}}$	Matthews correlation coefficient (MCC) $= \sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}} - \sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}}$	Threat score (TS), critical success index (CSI), Jaccard index $= \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$

Expected Utility

- Let y_x be the true label for x , and $P(y|x)$ the probability for label y our classifier predicts.
- Utility function: $U(y|y_x) \in \mathbb{R}$ [expresses how much we would value the prediction y if the true label is y_x]

What could be the goal now?

- Expected utility: $\mathbb{E}[U|y_x] = \sum_y U(y|y_x) \cdot P(y|x)$
- Goal: Find a classifier that maximises $\sum_x \mathbb{E}[U|y_x]$
- Example:
$$U(y|y_x) = \begin{cases} -100, & \text{if } y \neq y_x \text{ and } x = (\text{age} > 60) \\ 0, & \text{otherwise} \end{cases}$$

→ Min. Errors for over 60-year-olds

Quality Measures for Regression

How do we measure quality of our predictions?



Common Error Metrics

Target labels $y_1, \dots, y_n \in \mathbb{R}$, predicted labels $\hat{y}_1, \dots, \hat{y}_n \in \mathbb{R}$

- **Mean Absolute Error (MAE):** $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| = \|y - \hat{y}\|_1$
The average (absolute) error between predicted and actual label.
Problematic: \rightarrow Relative size dependent on the domain.
 \rightarrow No differentiation between one huge error and many small ones.
- **Root Mean Squared Error (RMSE):** $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \|y - \hat{y}\|_2$
Properties: $\rightarrow \text{MAE} \leq \text{RMSE} \leq \frac{1}{\sqrt{n}} \text{MAE}$.
 \rightarrow Punishes larger errors more than MAE.
 \rightarrow High variance in errors: $\text{MAE} \ll \text{RMSE} \rightarrow$ can indicate outliers

Let's take a Quiz...

... go to Mentimeter!

Let's take a break...

Back on in 5 min!



Source:
“The Ethical Algorithm”
(Chapter 2) by Kearns & Roth

Privacy

Features

Sensitive Features = Personal and not necessarily desired to be published in connection with identifiers
... or recoverable (**re-identification**)!

Examples: grades, ethnicity, health issues, sexual orientation, ...

Explicit Identifiers = Allow identifying a person without other infos

Examples: name, birth number, address, ...

Deleting Explicit Identifiers ... is Not Enough!

Group Insurance Commission (GIC)
Released data in mid-90's for academic research:

- Records of hospital visits of state employees in Massachusetts
- Included: Birthdates, ZIP code, sex
- Removed: Explicit patient identifiers (names, addresses, ...)

Latanya Sweeney: Identified then Massachusetts governor!

- Purchased some voter data for 20\$: Names, address, sex
- 6 people with same birth date, 3 of men, 1 correct ZIP code

... and showed: 87 % of Americans can be uniquely identified by birth date, ZIP code & sex!





Deleting User Information... is Not Enough!

The Netflix Price Competition (2006)

1 Mil \$ Prize for 10% improvement of recommendation accuracy based on released data:

- > 100mil user ratings of movies (1-5), ~0.5mil users, 18k movies
- Included: Generic user ID, user ratings of movies, time of rating
- Removed: Any other user information (no demographics!)

Arvind Narayanan, Vitaly Shmatikov:

Could identify (a small sets of matching) user records of an individual (after only 2 weeks)!

- Even if: Only 6 ratings are known with only approximate dates (± 14 days) \rightarrow 99% of records uniquely identified
- This is very possible by cross-referencing with public data (e.g., IMDB) ... on a large scale!

... Is this a privacy violation?!

\rightarrow Yes: Users did not publish *all* ratings.

K - Anonymity



- Idea: *Delete or bucket* data such that no record is unique!
- **k -anonymity**: For any query of features, e.g., ($25 < \text{Age} \leq 30$, Male, Norwegian) there are at least k records with the same features.
- Example: Making data 2-anonymous

Name	Post Code	Age	Gender	Nationality	Condition
*	*	$20 < \text{Age} \leq 25$	Female	English	Asthma
*	*	$25 < \text{Age} \leq 30$	Female	Norwegian	Diabetes
*	*	$25 < \text{Age} \leq 30$	Male	Norwegian	Chlamydia
*	*	$20 < \text{Age} \leq 25$	Female	English	Diabetes
*	*	$25 < \text{Age} \leq 30$	Female	Norwegian	Chlamydia
*	*	$25 < \text{Age} \leq 30$	Male	Norwegian	HIV
*	*	$25 < \text{Age} \leq 30$	Male	Norwegian	Chlamydia

K - Anonymity



- Idea: *Delete or bucket* data such that no record is unique!
- **k -anonymity**: For any query of features, e.g., $(25 < \text{Age} \leq 30, \text{Male}, \text{Norwegian})$ there are at least k records with the same features.
- Problem: Can still find out *some* information.
If Einar is $(25 < \text{Age} \leq 30, \text{Male}, \text{Norwegian}) \rightarrow$ He has Chlamydia *or* HIV!

Name	Post Code	Age	Gender	Nationality	Condition
*	*	$20 < \text{Age} \leq 25$	Female	English	Asthma
*	*	$25 < \text{Age} \leq 30$	Female	Norwegian	Diabetes
*	*	$25 < \text{Age} \leq 30$	Male	Norwegian	Chlamydia
*	*	$20 < \text{Age} \leq 25$	Female	English	Diabetes
*	*	$25 < \text{Age} \leq 30$	Female	Norwegian	Chlamydia
*	*	$25 < \text{Age} \leq 30$	Male	Norwegian	HIV
*	*	$25 < \text{Age} \leq 30$	Male	Norwegian	Chlamydia

K - Anonymity



- Idea: *Delete or bucket* data such that no record is unique!
- **k -anonymity**: For any query of features, e.g., (male, 24, Norwegian) there are at least k records with the same features.
- Problem: Multiple k -anonymous data might not be k -anonymous together!
If Einar is (25 < Age ≤ 30, Male, Norwegian, Student) → He has Chlamydia!

Age	Profession	Condition
20 < Age ≤ 25	Student	Asthma
25 < Age ≤ 30	Student	Diabetes
25 < Age ≤ 30	Student	Chlamydia
20 < Age ≤ 25	Student	Diabetes
25 < Age ≤ 30	Student	Chlamydia
25 < Age ≤ 30	Nurse	HIV
25 < Age ≤ 30	Nurse	Chlamydia

Age	Gender	Nationality	Condition
20 < Age ≤ 25	Female	English	Asthma
25 < Age ≤ 30	Female	Norwegian	Diabetes
25 < Age ≤ 30	Male	Norwegian	Chlamydia
20 < Age ≤ 25	Female	English	Diabetes
25 < Age ≤ 30	Female	Norwegian	Chlamydia
25 < Age ≤ 30	Male	Norwegian	HIV
25 < Age ≤ 30	Male	Norwegian	Chlamydia

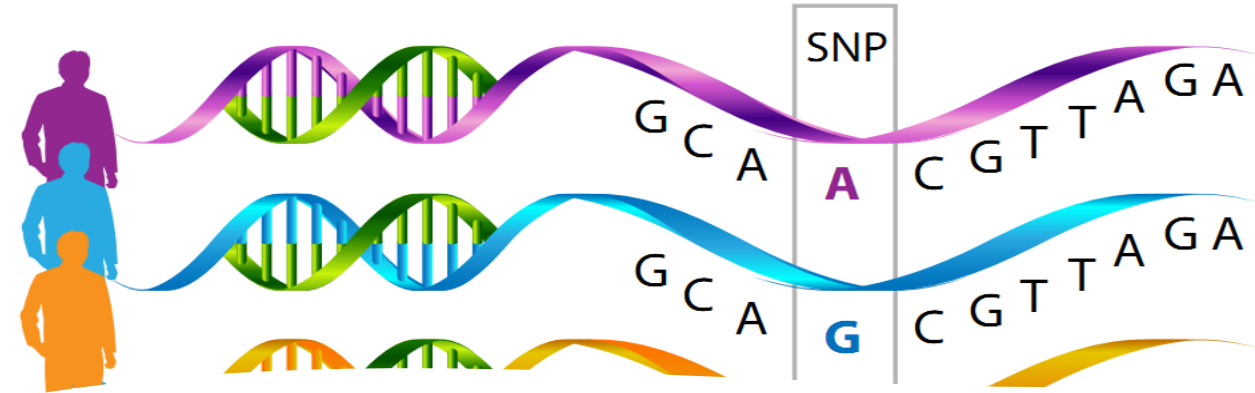
Multiple Data Sources

- Making data anonymous is not enough
 - One might cross-reference with other data available
 - We cannot anticipate other / future data resources!
- What can we do?
 - Only publishing statistics?
 - Adding random noise?

Publishing Statistics

Genome Wide Association Study (GWAS):

- Find correlations between genetic variances and prevalence of a disease
- Collect disease-patient's DNA data: genomes differ in only few positions
- Publish frequencies (of variances in genome positions) across population → Many such averages!

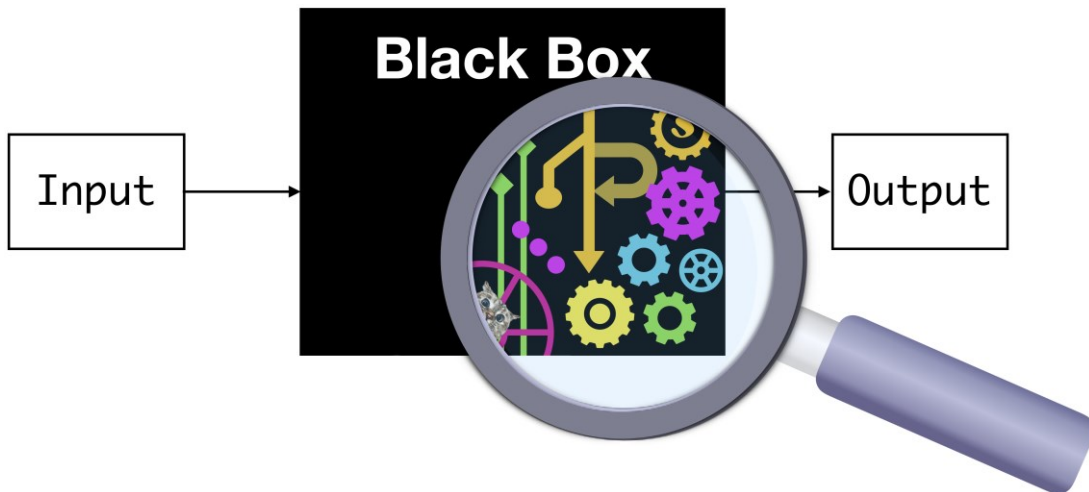


Homer, Szelling et al. in PLOS Genetics (2008):

Can test whether a particular person participates in a GWAS!

→ Privacy concern: GWAS participants have some disease!

→ Such data no longer open access (limits for research!)



Publishing Trained Models

Given a trained ML model (input – output observations):

→ Often possible to identify data points the model was trained on

→ Training data often has better prediction than non-training data (overfitting)

Differential Privacy

R. Doll, A. B. Hill study: Smoking and lung cancer

- 2/3 of all registered physicians in the UK in 1951 participated
- 1956: Strong evidence that smoking increases risk of lung cancer
- Participant (Roger, Physician, smoker) at increased risk of lung cancer
 - Harmful information (increased insurance costs, ...)
- But: This information would be harmful whether they participated or not!
 - No privacy violation!

Differential Privacy (DP): A study is DP if what can be learned about a person (or the harm they get) is (almost) the same whether or not they are included.





ϵ - Differential Privacy

- Randomized algorithm \mathcal{A} :

$$D_1 = (d_1, \dots, d_{i-1}, d_i, d_{i+1}, \dots, d_n)$$

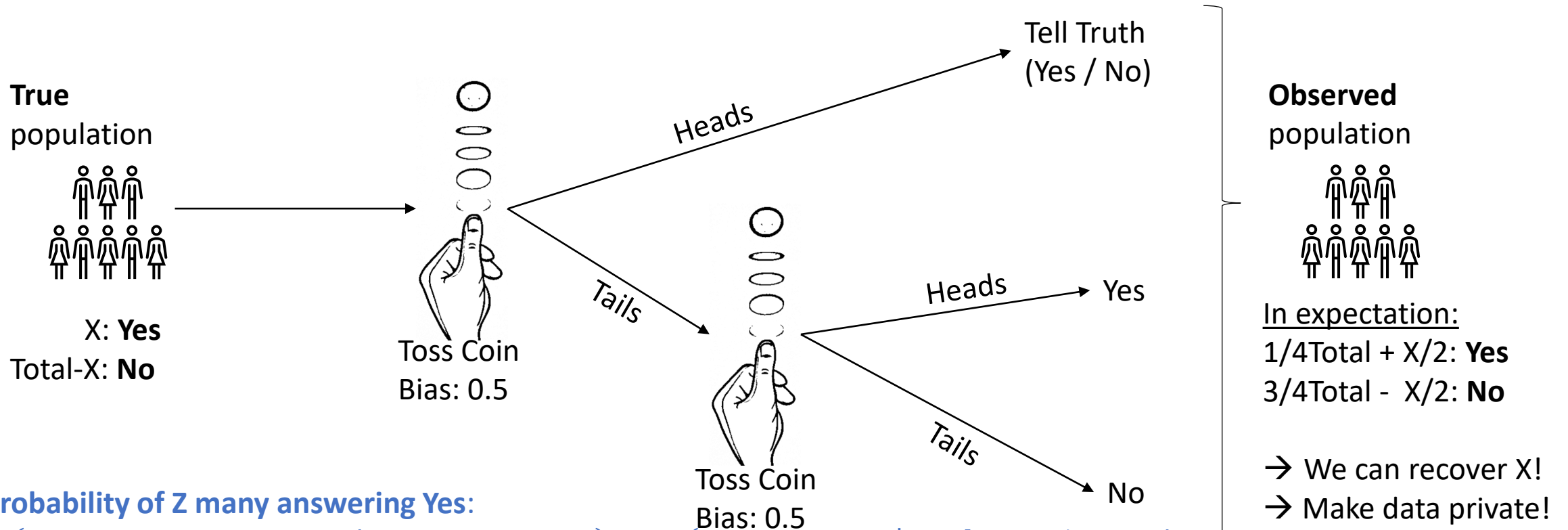
$$\parallel \quad \dots \quad \parallel \quad \# \quad \parallel \quad \dots \quad \parallel$$

$$D_2 = (d_1, \dots, d_{i-1}, \hat{d}_i, d_{i+1}, \dots, d_n)$$



$$\rightarrow \mathcal{A} \text{ is } \epsilon\text{-differential private if: } P[\mathcal{A}(D_1) \in S] \leq e^\epsilon \cdot P[\mathcal{A}(D_2) \in S]$$
$$\Leftrightarrow \frac{P[\mathcal{A}(D_1) \in S]}{P[\mathcal{A}(D_2) \in S]} \leq e^\epsilon$$

Randomized Response Mechanism



Probability of Z many answering Yes:

$$\frac{P(Z \text{ many Yes-answers} \mid X \text{ many Yes})}{P(Z \text{ many Yes-answers} \mid X-1 \text{ many Yes})} = \frac{P(\text{answer=Yes} \mid \text{truth=Yes})}{P(\text{answer=Yes} \mid \text{truth=No})} = \frac{3/4}{1/4} = 3 = e^\epsilon$$

Probability of Z many answering Yes:

$$\frac{P(Z \text{ many Yes-answers} \mid X \text{ many Yes})}{P(Z \text{ many Yes-answers} \mid X+1 \text{ many Yes})} = \frac{P(\text{answer=Yes} \mid \text{truth=Yes})}{P(\text{answer=No} \mid \text{truth=Yes})} = \frac{3/4}{1/4} = 3 = e^\epsilon$$

$\Rightarrow \epsilon = \ln(3)$ -differential privacy

$\Rightarrow \epsilon = \ln(3)$ -differential privacy

Summary

Measuring Quality of Classification

Predicted Class	Actual Class	
	+	-
	True Positives (TP)	False Positives (FP)
	False Negatives (FN)	True Negatives (TN)

Accuracy:

$$\frac{TP + TN}{Total}$$

Precision:

$$\frac{TP}{TP + FP}$$

Recall:

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

Measuring Quality of Regression

Mean Absolute Error (MAE):

$$\frac{1}{n} \sum_{i=0}^{n-1} |y_i - \hat{y}_i| = \|y - \hat{y}\|_1$$

The average (absolute) error between predicted and actual label.

Root Mean Squared Error (RMSE):

$$\sqrt{\frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2} = \|y - \hat{y}\|_2$$

$$MAE \leq RMSE \leq \frac{1}{\sqrt{2}} MAE.$$

Privacy

Deleting identifiers:

- Not obvious how much / what to delete
- Cross-referencing with other datasets can reveal identities

k-anonymity:

Every row same features as k others.

- Cross-referencing with other datasets can reveal identities

ε-differential privacy:

Whether one person is included or not (or changes response or not) makes little difference to what we can learn.

→ Randomized response mechanism:
Make data $\ln(3)$ -diff. private by adding noise with help of coin flips!

Which measure fits best? – Application dependent!
To grasp the quality, look at several!