

INFO 251: Applied Machine Learning

# Missing and Imbalanced Data

# Announcements

- Problem set 6 due April 17
- Questions about last two lectures?
  - Kent is graciously willing to take follow-up questions:  
[kentkchang@berkeley.edu](mailto:kentkchang@berkeley.edu)

# Course Outline

- Causal Inference and Research Design
  - Experimental methods
  - Non-experiment methods
- **Machine Learning**
  - Design of Machine Learning Experiments
  - Linear Models and Gradient Descent
  - Non-linear models
  - Fairness and Bias in ML
  - Neural models
  - Deep Learning
  - **Practicalities**
  - Unsupervised Learning
- Special topics

# Key Concepts: LLMs

- Distributional hypothesis
- Word2Vec
- Embeddings
- Recurrent Neural Networks (RNNs)
- Long Short-Term Memory (LSTM)
- Attention and Self-Attention
- Transformers
- Model Context Protocol (MCP)
- Controlling LLMs

# Intuition check

- Which of the following models are good at encoding local and long-range dependencies? (For example, would be good at distinguishing between different meanings of “left”)
  - Convolutional Neural Networks (CNNs)
  - Word2Vec embeddings
  - Recurrent Neural Networks (RNNs)
  - Long Short-Term Memory (LSTM)
  - Attention

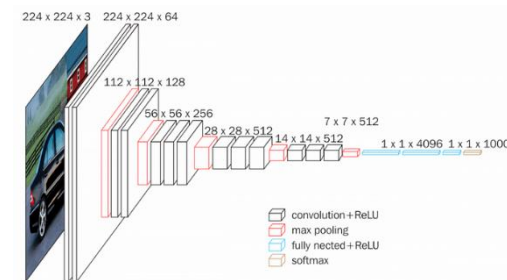
# Outline

- Deep Learning: Quick Recap
- Imbalanced data: 4 ideas
  - Simple things
  - Resampling
  - Weighting
  - Algorithm-level adjustments
- Missing data
  - Dropping
  - Imputing

# Background: CNNs, Embeddings

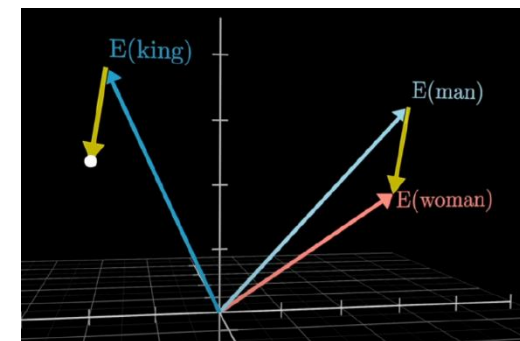
## ■ CNNs

- Input is fixed-sized vector (e.g., an image); output is a fixed-sized vector as output (e.g., class probabilities)
- Ordering of inputs is largely irrelevant; “no persistence” and “no memory”



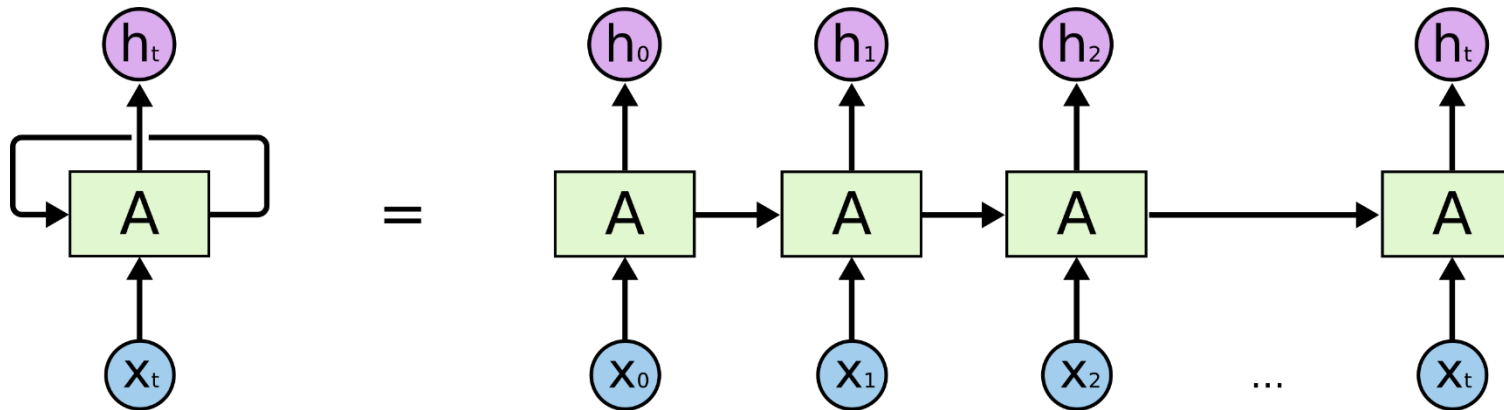
## ■ Embeddings

- A learned vector representation of input data (e.g., text)
- In NLP: Embeddings map tokens (like words) to vectors that encode their meaning in context-aware ways
- But embeddings are *static*, don't capture local context



# Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
  - A class of neural networks “with loops”, i.e., where previous outputs can be used as inputs

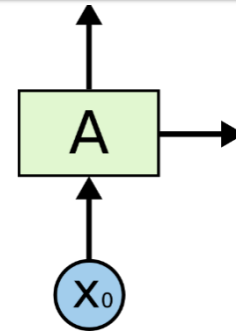


- Each node accepts input vector, emits output vector
  - Output influenced by input  $x$ , as well as prior state



# Recurrent Neural Networks

- Computation: Unpacking the

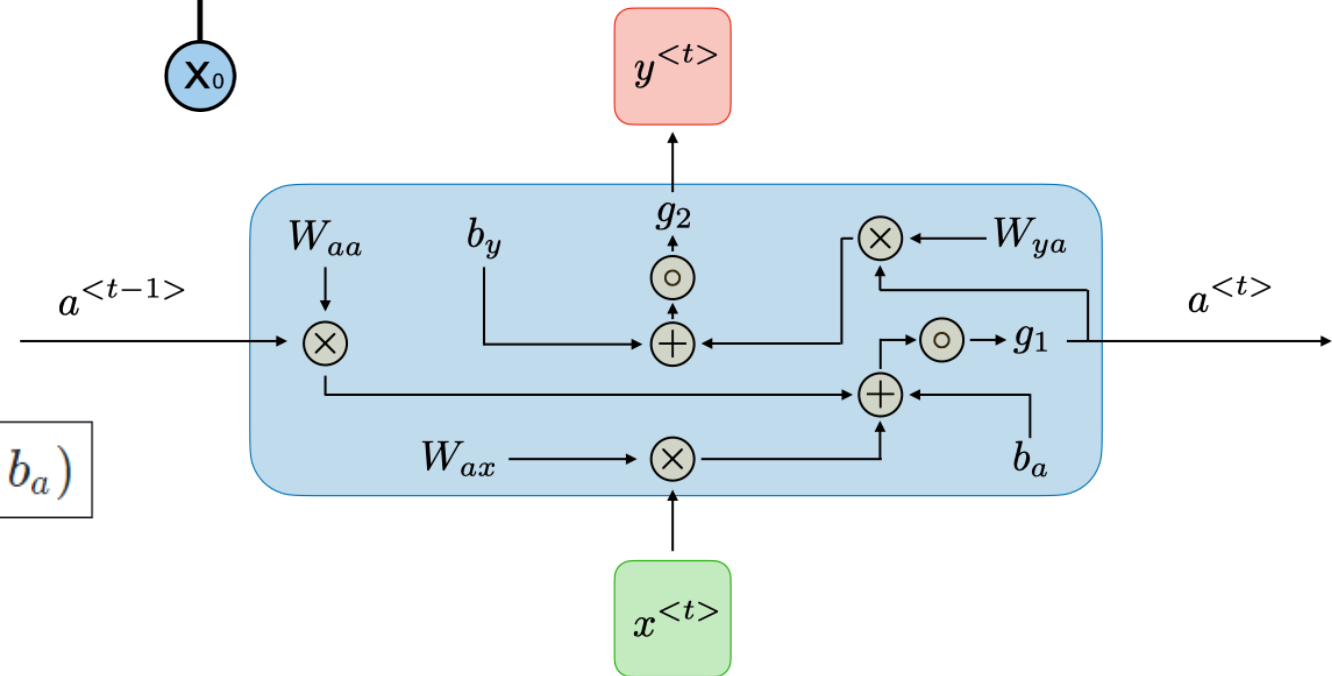


- For each timestep  $t$ :

$$a^{<t>} = g_1(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

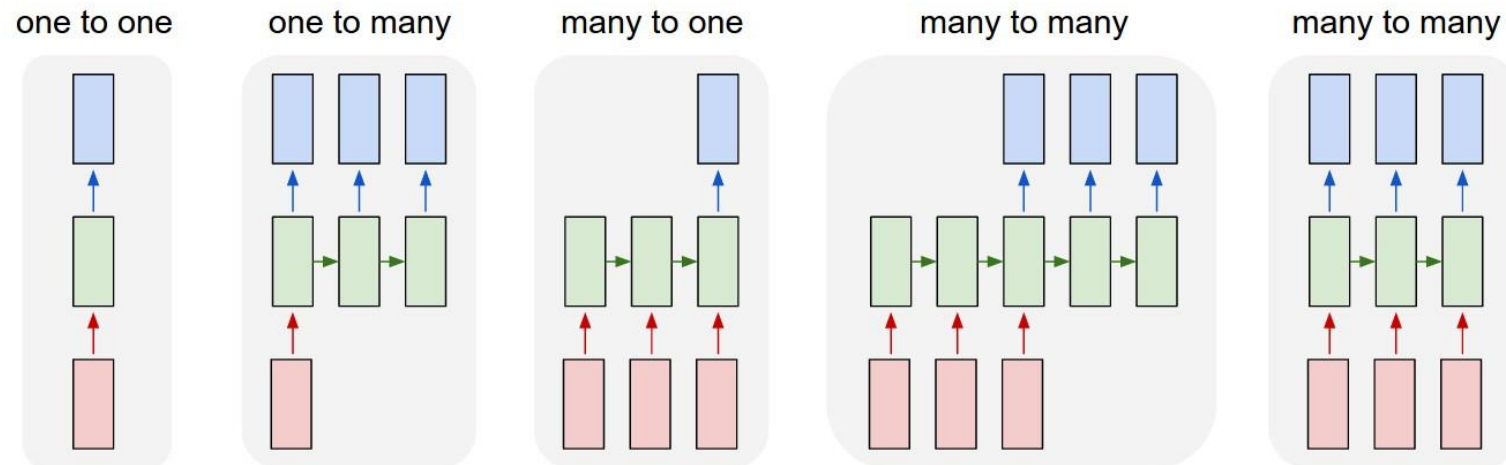
$$y^{<t>} = g_2(W_{ya}a^{<t>} + b_y)$$

where  $W_{ax}, W_{aa}, W_{ya}, b_a, b_y$  are coefficients that are shared temporally  
 $g_1, g_2$  activation functions.



# Example RNN architectures

- Recurrent Neural Networks (RNNs)
  - Current classification depends on past information
  - Makes it possible to operate over *sequences* of vectors: Sequences in input, output, or both



# From RNN's to LSTM

## ■ Some pros and cons of RNN's

Advantages	Drawbacks
<ul style="list-style-type: none"> <li>• Possibility of processing input of any length</li> <li>• Model size not increasing with size of input</li> <li>• Computation takes into account historical information</li> <li>• Weights are shared across time</li> </ul>	<ul style="list-style-type: none"> <li>• Computation being slow</li> <li>• Difficulty of accessing information from a long time ago</li> <li>• Cannot consider any future input for the current state</li> </ul>

<https://cs230.stanford.edu/>

## LONG SHORT-TERM MEMORY

NEURAL COMPUTATION 9(8):1735–1780, 1997

Sepp Hochreiter

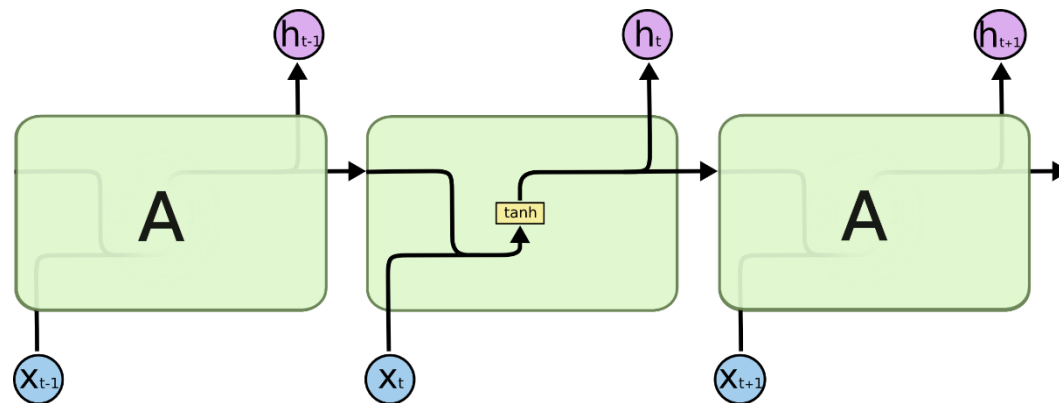
Jürgen Schmidhuber

## ■ Enter the LSTM

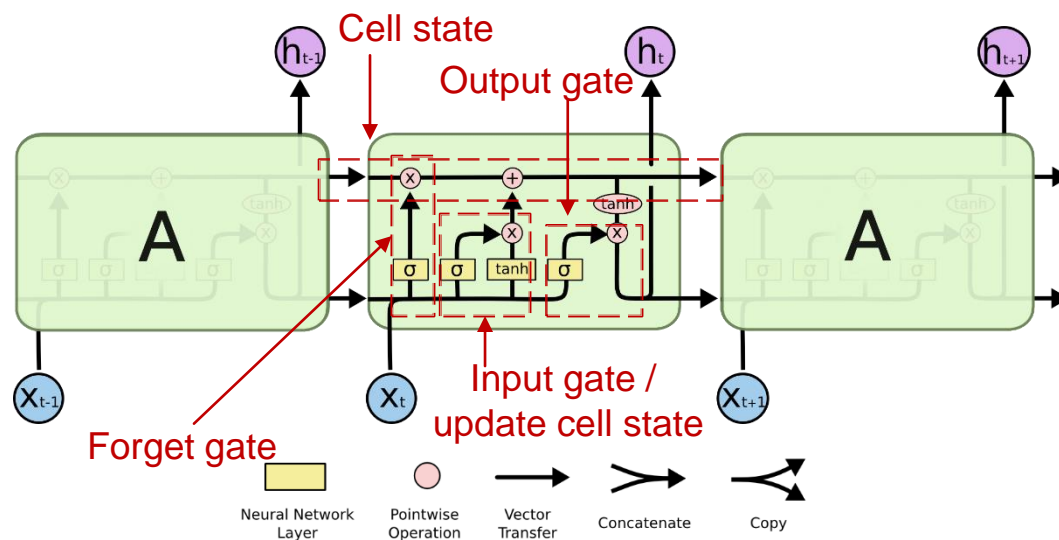
- *"An LSTM layer consists of a set of recurrently connected blocks, known as memory blocks. Each one contains one or more recurrently connected memory cells and three multiplicative units – the input, output and forget gates – that provide continuous analogues of write, read and reset operations for the cells."*

# LSTM architecture

- Standard RNN



- Standard LSTM

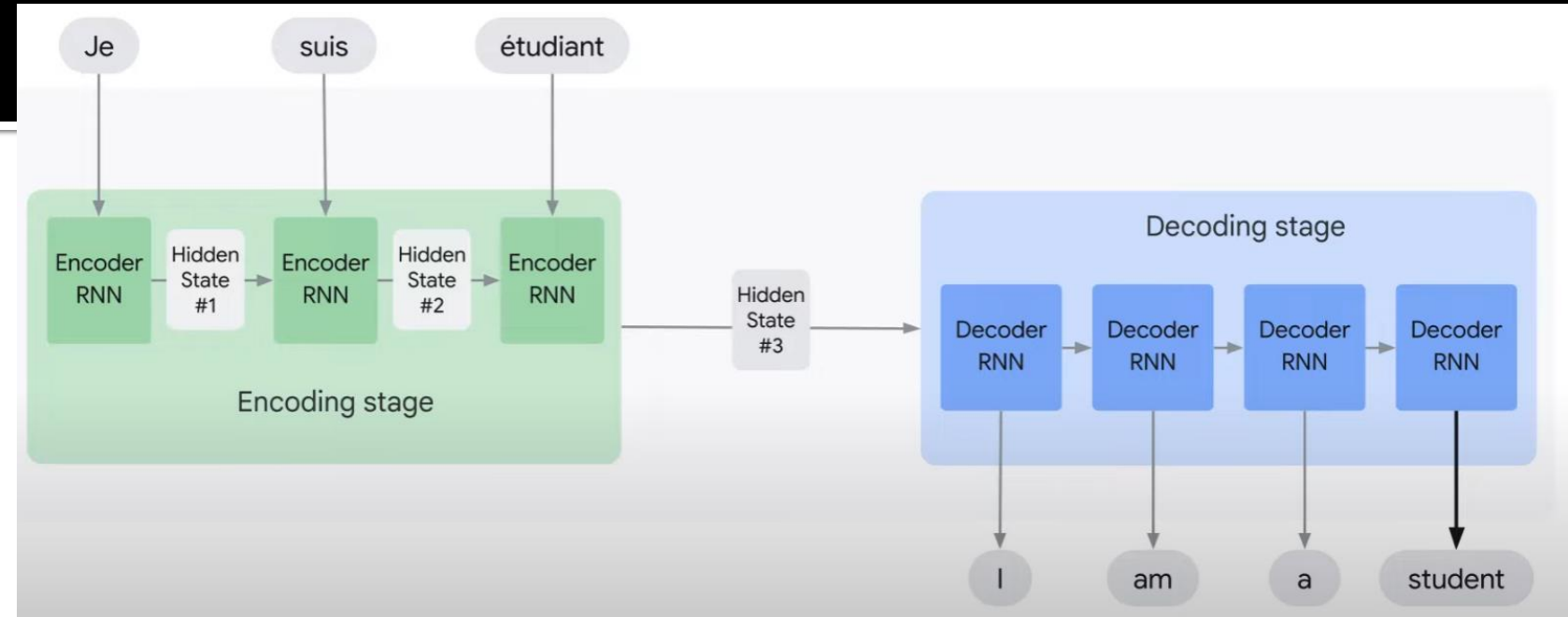


# From LSTM's to Attention

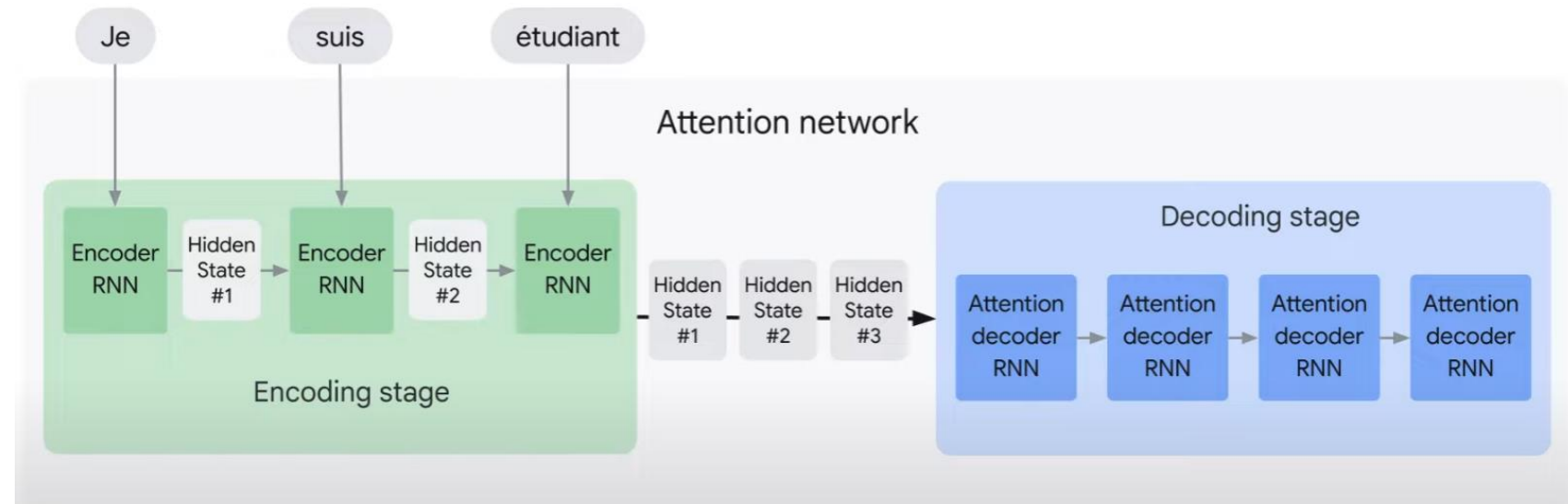
- Limitations of LSTM's
  - **Encoding bottleneck:** Single hidden state vector for context-dependencies. For long sequences, not all information fits well into that vector
  - **Sequential operations:** LSTMs process step-by-step, hard to parallelize
  - **No long-range dependencies:** Over longer distances, LSTMs can struggle to retain relevant information
- Attention: **Encodes context**
  - Allows model to look at **all encoder hidden states** and assign a **learned weight** to each one based on how relevant it is to the current decoding step
  - In other words, model “attends” to all/different hidden states

# Attention

Traditional RNN encoder-decoder



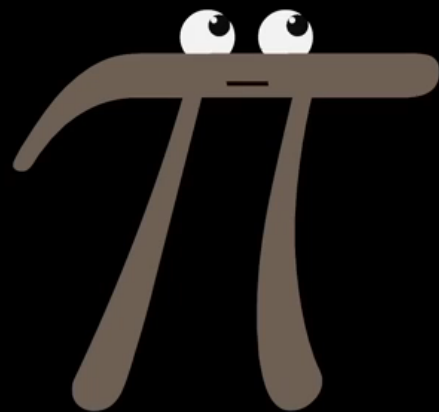
Attention model



# Summary

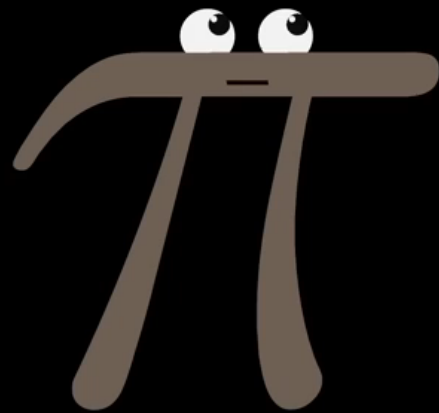
Model / Innovation	Key Concepts	Limitations Addressed
RNN	Sequential processing	Vanishing gradients, limited memory
LSTM	Gating mechanisms	Longer memory, but still sequential and bottlenecked
Attention	Learn to focus on parts of input	Removes fixed-size context bottleneck
Self-Attention	Each token attends to all others	Captures global context, fully parallel
Transformer	Stacks of self-attention + FFN	No recurrence, fast, scalable

# Transformers and attention





# Tuning



# Further Reading

## ■ Resources on bCourses

- Dive Into Deep Learning, Chapters 9-10 (RNN's) and 11 (Attention/Transformers)
- ConvNets for Visual Recognition: <http://cs231n.github.io/>
- UVA tutorial on transformers and attention: <https://uvadlc.github.io/> (tutorial week #5)

## ■ At Berkeley

- STAT 157: Introduction to Deep Learning
- CS 285: Deep Reinforcement Learning
- CS 294-131: Special Topics in Deep Learning
- CS294-158: Deep Unsupervised Learning

# Outline

- Deep Learning: Recap
- **Imbalanced data: 4 ideas**
  - Simple things
  - Resampling
  - Weighting
  - Algorithm-level adjustments
- Missing data
  - Dropping
  - Imputing

# Key concepts: Missing and Imbalanced Data

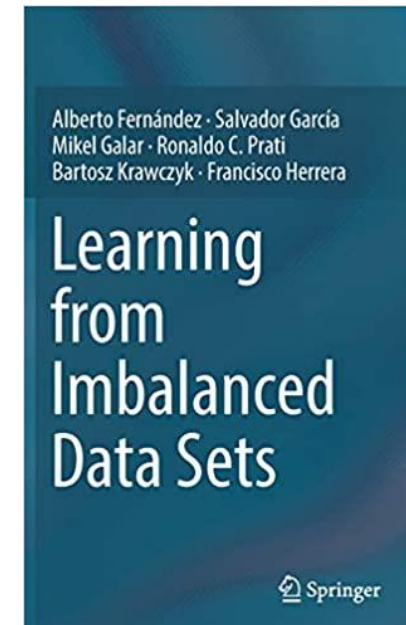
- Stratified randomization
- Upsampling and downsampling
- SMOTE and AdaSyn
- Reweighting
- Algorithm-level adjustments for imbalance
- Problems with data missingness
- Selective labels
- Model-free imputation (e.g., zero-coding, mean imputation)
- Model-based imputation (e.g., hot deck)

# Outline

- Deep Learning: Recap
- **Imbalanced data: 4 ideas**
  - Simple things
  - Resampling
  - Weighting
  - Algorithm-level adjustments
- Missing data
  - Dropping
  - Imputing

# Imbalanced data

- In many real-world instances, rate of true positives and true negatives are far from even
  - Credit card fraud
  - Disease diagnosis
  - Product adoption and churn
  - Terrorist threats
- *How to learn from imbalanced data?*



# Imbalanced data

- Note that many learning algorithms will fail if used “out of the box” on imbalanced data
  - Most algorithms minimize *error*
  - A fast path to error minimization is predicting the majority class
    - Daume: *If a teacher told you to study for an exam with 1000 True/False questions and you knew that only one question had a correct answer of True, how much would you study?*

# Imbalanced data

- Idea 1: Simple things
  - Don't rely on accuracy as a performance metric
  - Instead, use
    - Confusion matrices
    - ROC curves
    - Cohen's Kappa (normalizes accuracy by imbalance)

$$k = \frac{p_o - p_e}{1 - p_e}$$

- $p_o$  = actual accuracy
- $p_e$  = chance accuracy

- F-scores, precision, recall, etc.

		Predicted:		
		NO	YES	
Actual:	NO	TN = 50	FP = 10	60
	YES	FN = 5	TP = 100	105
		55	110	



# Imbalanced data

- Idea 1: Simple things
  - Don't rely on accuracy as a performance metric
  - **Use appropriate baselines**

	Accuracy	Recall	Precision	F	AUC	% Answered Yes
<i>Panel A: Assets and Housing</i>						
Owens a radio	0.976	1.000	0.976	0.988	0.899	0.973
Owens a bicycle	0.676	0.552	0.678	0.609	0.722	0.456
Household has electricity	0.819	0.533	0.761	0.627	0.828	0.285
Owens a television	0.855	0.497	0.738	0.594	0.814	0.214
Has indoor plumbing	0.887	0.250	0.842	0.386	0.843	0.142
Owens a motorcycle/scooter	0.899	0.011	1.000	0.022	0.772	0.102
Owens a car/truck	0.945	0.213	0.867	0.342	0.849	0.068
Owens a refrigerator	0.954	0.180	1.000	0.305	0.878	0.055
Has landline telephone	0.992	0.125	1.000	0.222	0.562	0.009
<i>Panel B: Social Welfare Indicators</i>						
Hospital bills in last 12 months	0.633	0.890	0.633	0.740	0.653	0.587
Very ill in last 12 months	0.686	0.188	0.550	0.280	0.671	0.325
Death in family in last 12 months	0.665	0.183	0.632	0.284	0.619	0.363
Flood or drought in last 12 months	0.788	0.086	0.607	0.151	0.706	0.219
Fired in last 12 months	0.901	0.022	1.000	0.043	0.731	0.101

Table 2: Model performance at predicting responses from survey respondents based on call records data

# Imbalanced data

- Idea 1: Simple things
  - Don't rely on accuracy as a performance metric
  - Use appropriate baselines
  - **Stratified randomization**
    - When randomizing into test-train, separately randomize minority class observations and majority class observations – ensures even proportions in test/train
    - Same goes for cross-validation, when randomly assigning observations to folds

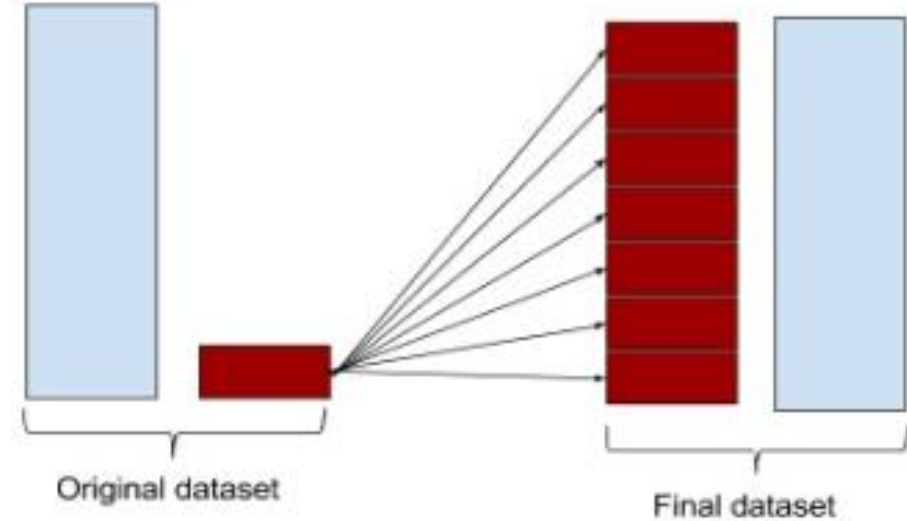
# Imbalanced data

- Idea 1: Simple things
  - Don't rely on accuracy as a performance metric
  - Use appropriate baselines
  - Stratified randomization
  - **Adjust classification threshold**
    - E.g. Provost & Fawcett: based on true class distributions and cost of mis-classifications
  - **Gather more data!**

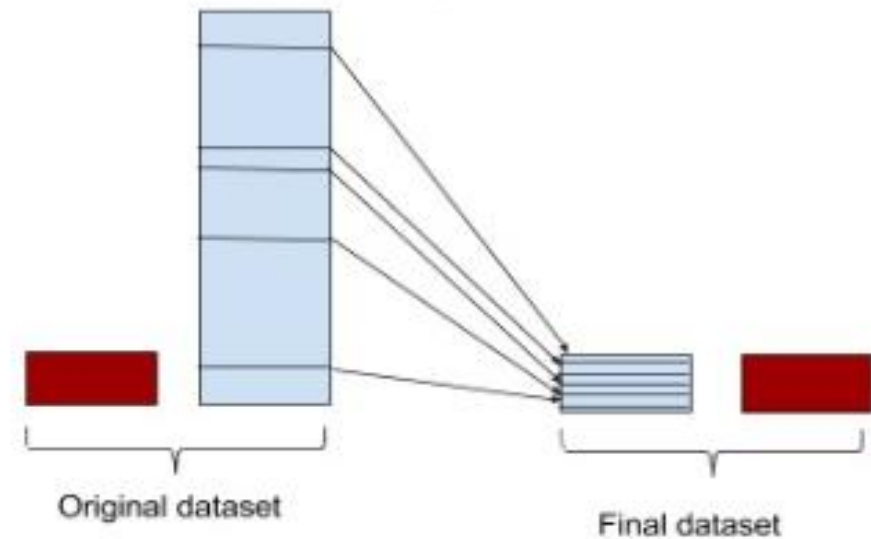
# Imbalanced data

- Idea 2: **Resampling**
  - Subsampling
  - Oversampling
  - Mixtures

**Oversampling** minority class



**Undersampling** majority class



# Imbalanced data

## ■ Idea 2: Resampling

### ■ Sub-sampling, or “Down-sampling”

- Idea: Include all instances of minority class; include each instance of majority class with probability  $1/\alpha$
- This involves throwing out some data
- Typically, set  $\alpha$  = ratio of majority/minority class, to achieve 50-50 split of training data

# Imbalanced data

- Idea 2: **Resampling**

- Oversampling, or “**Up-sampling**”

- Idea: Include all instances of majority class, include each instance of minority class  $\alpha$  times
    - All data is used, some data is used several times

# Imbalanced data

- Idea 2: **Resampling – Other options**
  - SMOTE: Synthetic Minority Over-sampling Technique
    - Creates artificial data based on the feature space similarities between existing minority examples.
  - ADASYN: ADaptive SYNthetic Sampling Approach
    - Creates more instances near mis-classified (minority) instances

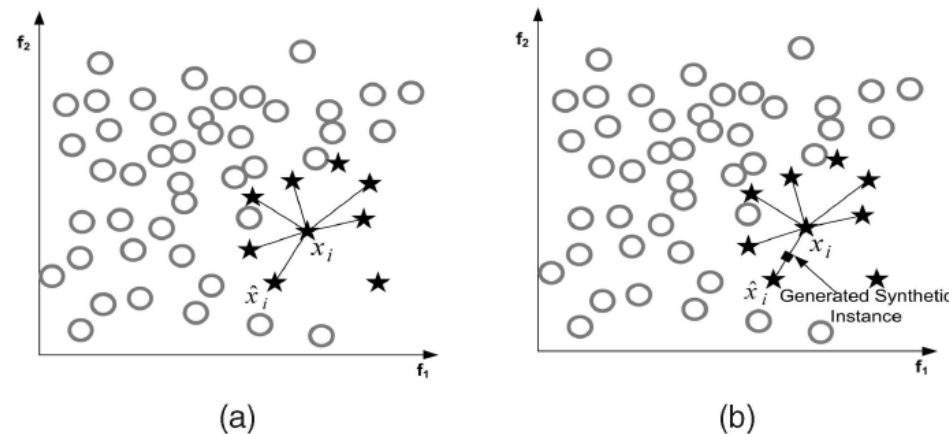


Fig. 3. (a) Example of the K-nearest neighbors for the  $x_i$  example under consideration ( $K = 6$ ). (b) Data creation based on euclidian distance.

# Imbalanced data

## ■ Idea 2: Resampling

### ■ Other considerations

- In both cases, *error rate* is ( $\alpha$  times) higher than it will be on unbalanced classifier -- but other measures of performance often improve
- Resampling typically done *after* train/test split

### ■ Comparing up- and down-sampling

- Up-sampling often produces lower error (it sees more variance in training data!)
- Computational efficiency is main reason for down-sampling
- You can try both!



# Imbalanced data

## ■ Idea 3: **Weighting**

- Instead of explicitly up- or down-sampling training data, associate a weight with each observation
- For many algorithms this is straightforward and computationally efficient
  - k-Nearest Neighbors
  - Decision Trees / Random Forests (remember Adaboost?)
  - Regression
- When possible, this is often the preferred approach

# Imbalanced data

- Idea 4: **Algorithm-level adjustments**
  - Modify learning algorithm to reduce bias towards majority class. For example:
    - “Cost-sensitive” approaches modify the learner to vary penalty for different classes of examples, different (mis-)classifications, etc.
- Many hybrid methods combine algorithm-level and data-level adjustments

# Imbalanced data

## ■ Further reading

- He, H., Garcia, E.A., 2009. Learning from Imbalanced Data. IEEE Transactions on Knowledge and Data Engineering 21, 1263–1284.
- Sun, Y., Wong, A.K.C., Kamel, M.S., 2009. Classification of imbalanced data: a review. Int. J. Patt. Recogn. Artif. Intell. 23, 687–719.

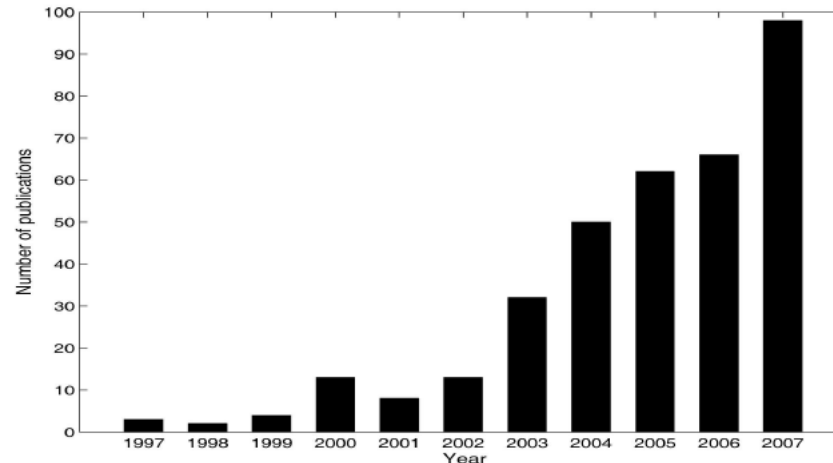


Fig. 1. Number of publications on imbalanced learning.

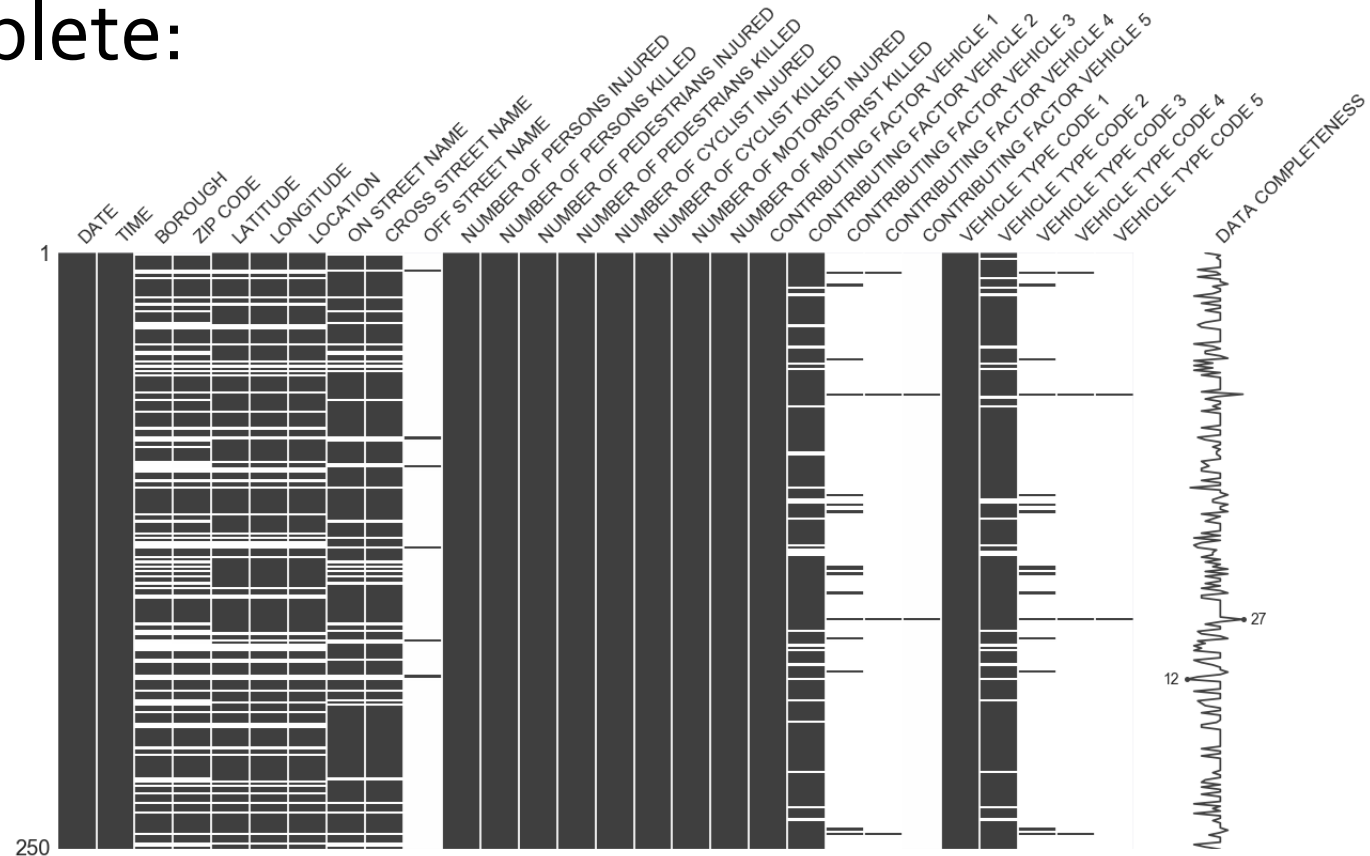
# Outline

- Deep Learning: Recap
- Imbalanced data: 4 ideas
  - Simple things
  - Resampling
  - Weighting
  - Algorithm-level adjustments
- **Missing data**
  - Dropping
  - Imputing
  - “Selective labels”

# Missing data

- A common issue in applied ML settings is that a data matrix is frequently not complete:

```
import missingno as msno
msno.matrix(myData)
```



# Missing data

- What to do?
  - Drop features
    - Great, if you can afford to do so!
  - Drop observations
    - Great, if you can afford to do so!
- In both cases, biggest concern is when data missingness is non-random
- Important: Many ML packages (e.g., sklearn) automatically drop observations and/or features with missing data!
  - Make sure to check!

# Systematic missingness

- Systematic missingness: Example
  - We are interested in predicting who will default on a loan. We drop all features where >50% of data are missing and train a predictive model:

$$default_i = f(X_i) + error_i$$

- Several features get dropped: `last_loan`, `last_loan_value`, `last_loan_APR`, `last_loan_repaid`
- What might go wrong here?

# Missing data

- What to do?
  - Drop features
    - Great, if you can afford to do so!
  - Drop observations
    - Great, if you can afford to do so!
  - Create new variables
    - E.g., `feature_K_is_present`, and `feature_K_is_present * feature_K`
  - Impute, replace



# Imputation

- Model-free approaches
  - Draw a value from distribution of  $X$
  - Zero-coding, filling with marker (-9999)
  - Mean imputation
  - Interpolation, “carry-forward” (for panel data)
- Model-based approaches (done with training data!)
  - Regression modeling
  - Matching (“hot deck imputation”)
  - k-NN imputation
  - (any other supervised algorithm can be used too!)

# Selective labels

- One final note: The “selective labels problem”
  - We are interested in predicting who will default on a loan:

$$\text{default}_i = f(X_i) + \text{error}_i$$

- We only can train on the population of people who have received loans in the past. Is this the same population for whom we want to make predictions?
- Common settings
  - Bail decisions, recidivism
  - Hiring decisions
  - Admissions decisions
  - Cf. Lakkaraju et al. (2017 KDD)

# Further reading

1. Introduction to KDD and Data Science
2. Foundations on Imbalanced Classification
3. Performance Measures
4. Cost-Sensitive Learning
5. Data Level Preprocessing Methods
6. Algorithm-Level Approaches
7. Ensemble Learning
8. Imbalanced Classification with Multiple Classes
9. Dimensionality Reduction for Imbalanced Learning
10. Data Intrinsic Characteristics
11. Learning from Imbalanced Data Streams
12. Non-classical Imbalanced Classification Problems
13. Imbalanced Classification for Big Data
14. Software and Libraries for Imbalanced Classification

