

Machine Learning

Lecture 7: Gradient boosting

Harbour.Space University
February 2020

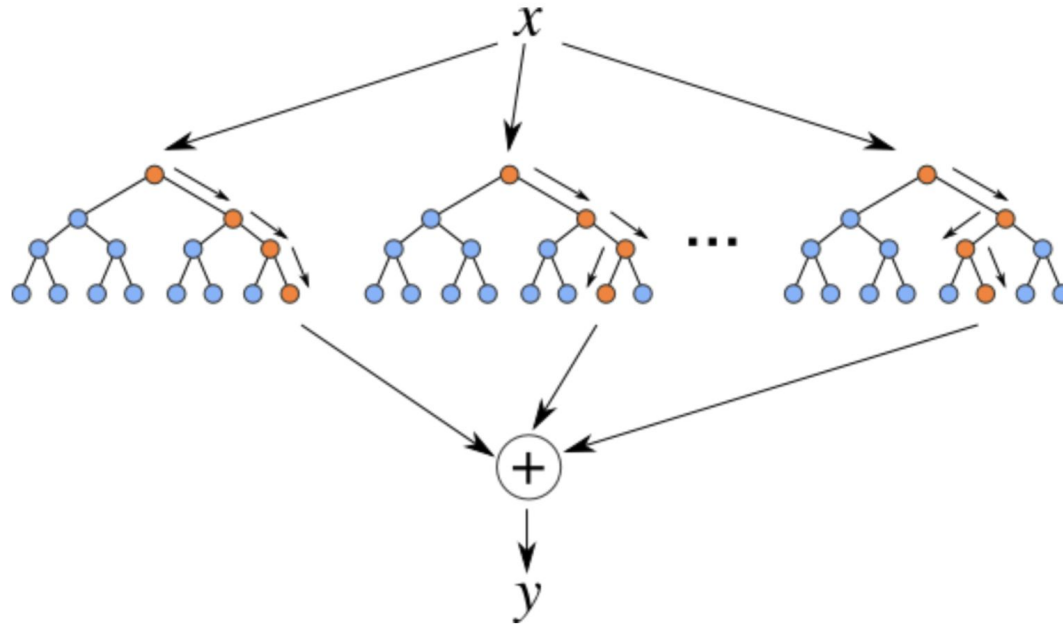
Radoslav Neychev

Outline

1. Boosting
2. Gradient boosting
3. Offtopic: word representations

Random Forest

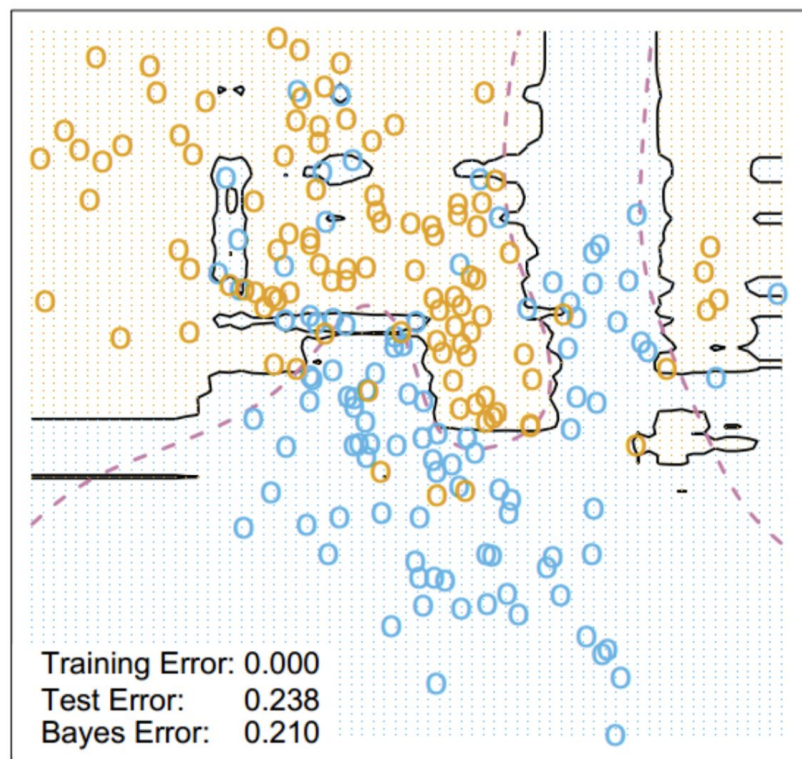
Bagging + RSM = Random Forest



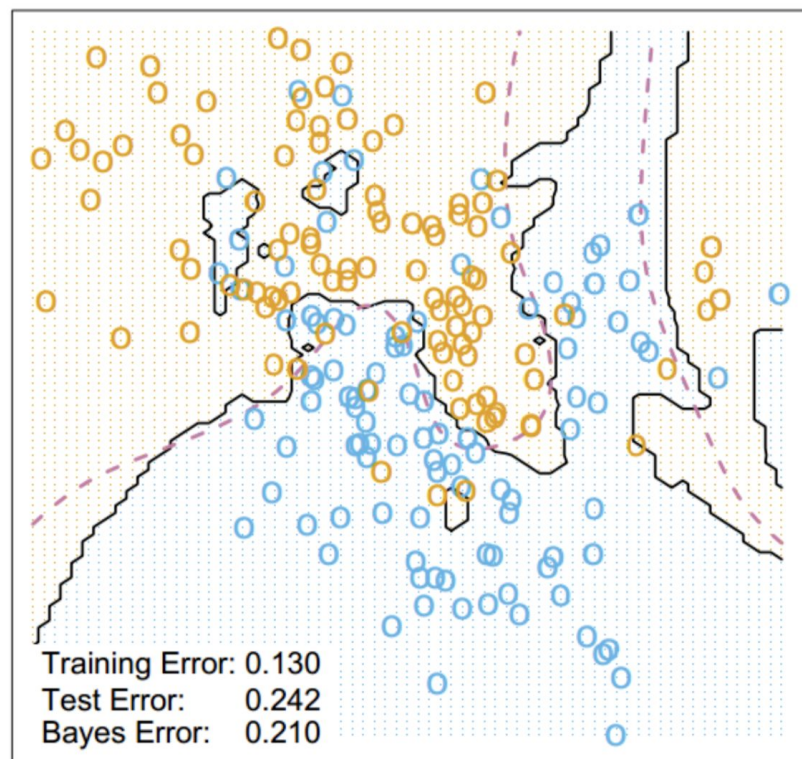
- One of the greatest “universal” models.
- There are some modifications: Extremely Randomized Trees, Isolation Forest, etc.
- Allows to use train data for validation: OOB

$$\text{OOB} = \sum_{i=1}^{\ell} L \left(y_i, \frac{1}{\sum_{n=1}^N [x_i \notin X_n]} \sum_{n=1}^N [x_i \notin X_n] b_n(x_i) \right)$$

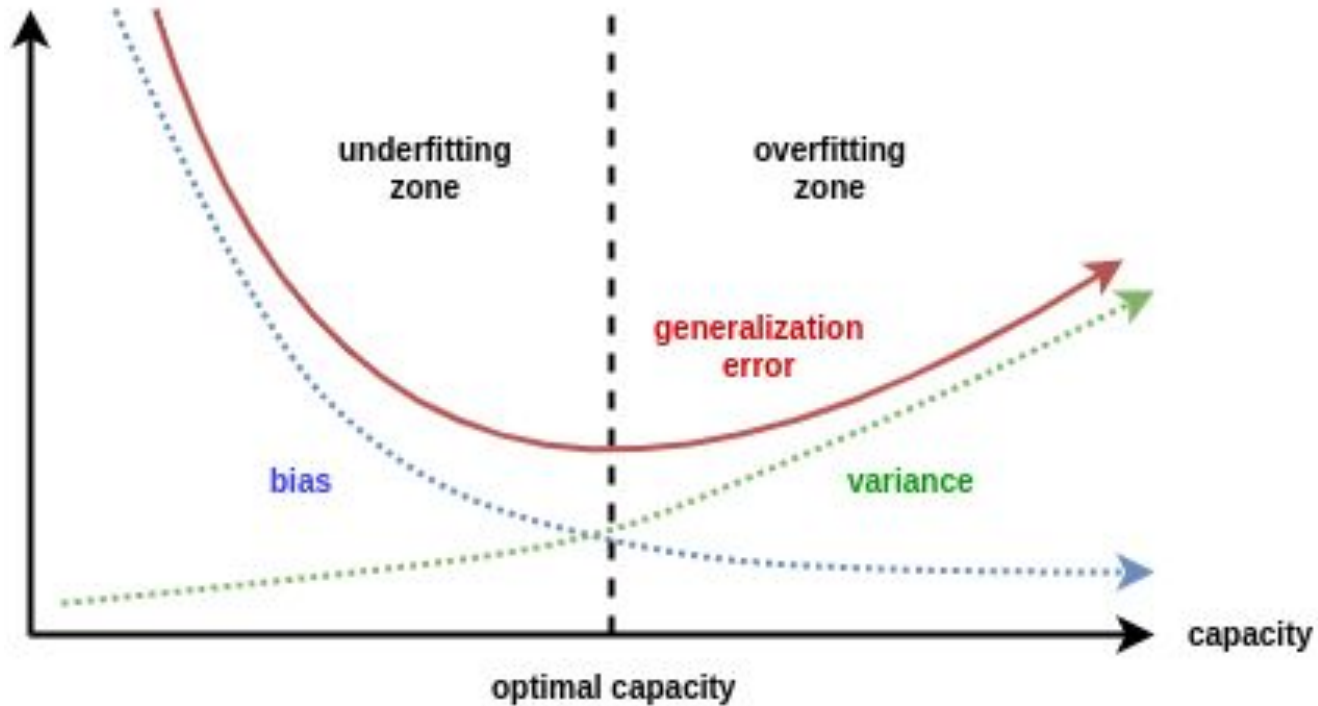
Random Forest Classifier



3-Nearest Neighbors



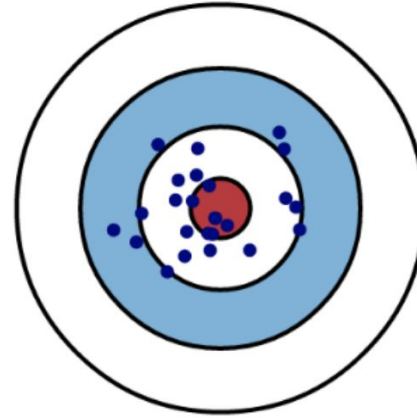
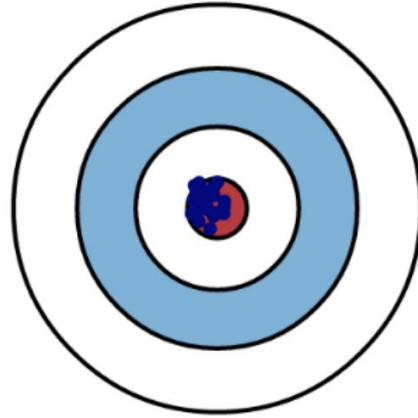
Bias-variance tradeoff



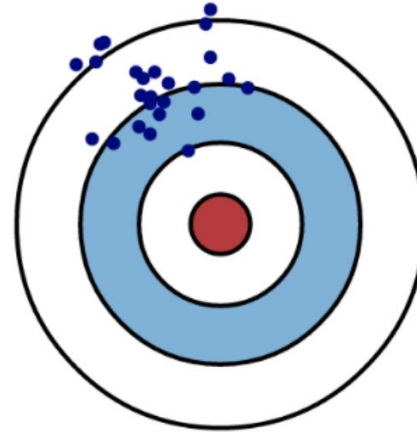
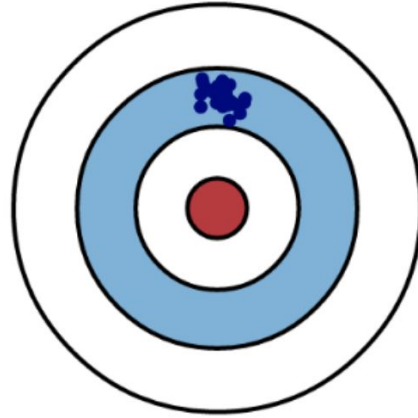
Low Variance

High Variance

Low Bias

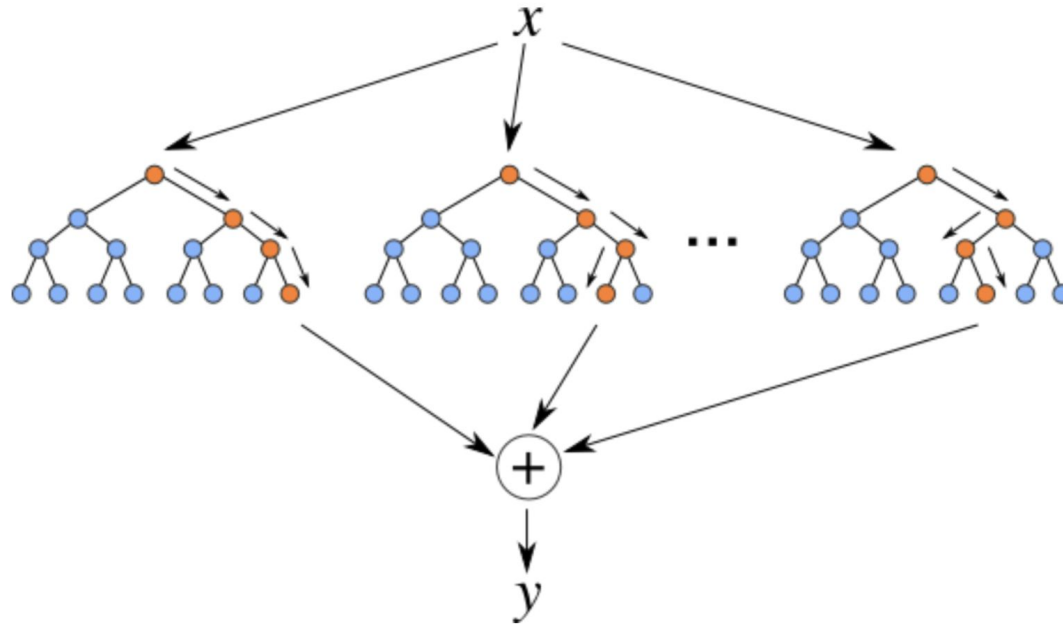


High Bias



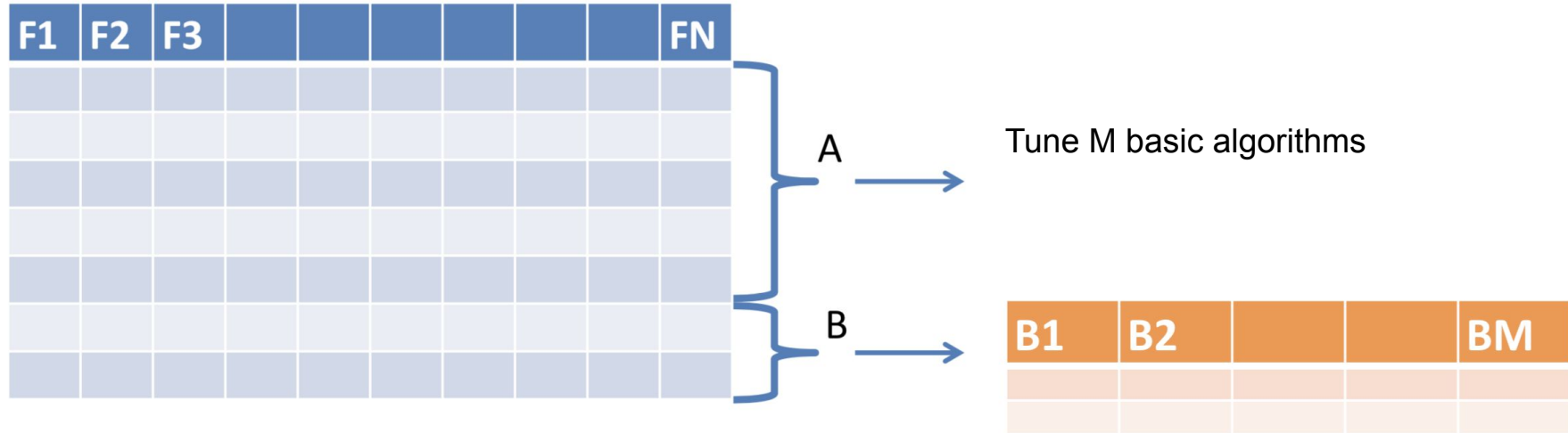
Random Forest

Is Random Forest decreasing bias or variance by building the trees ensemble?



Stacking

How to build an ensemble from *different* models?

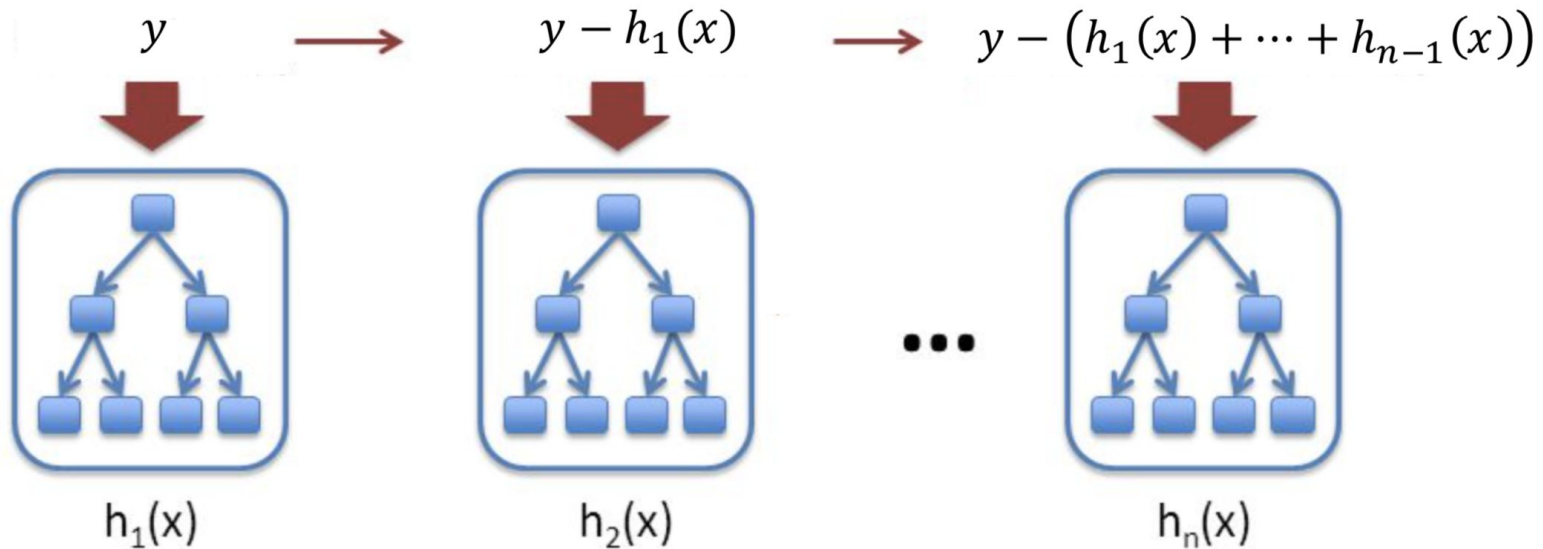


$$a(x) = \sum_{t=1}^T \alpha_t b_t(x)$$

e.g.

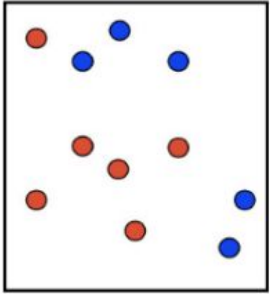
Gradient boosting

$$a_n(x) = h_1(x) + \cdots + h_n(x)$$

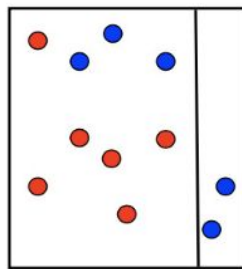
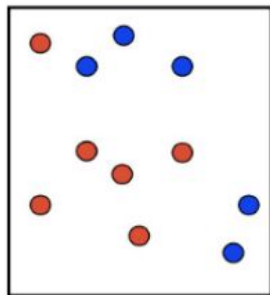


Boosting: intuition

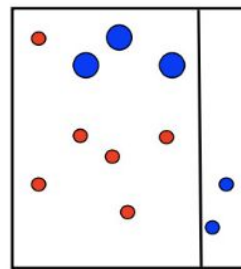
Binary classification problem.
Models - decision stumps.



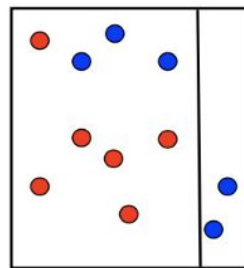
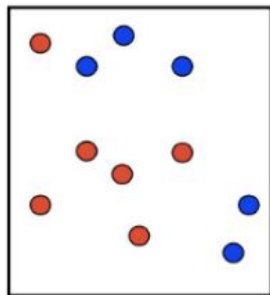
Boosting: intuition



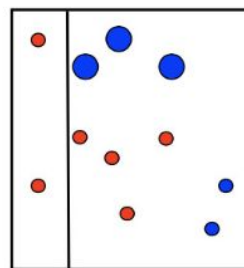
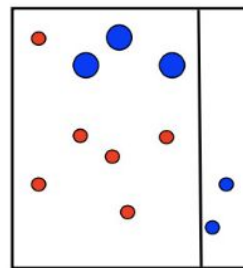
$t = 1$



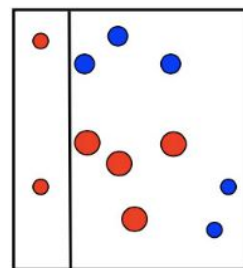
Boosting: intuition



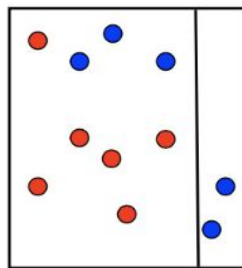
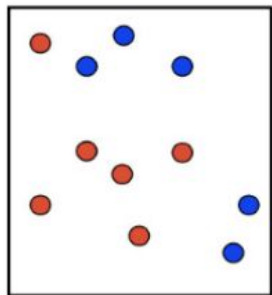
$t = 1$



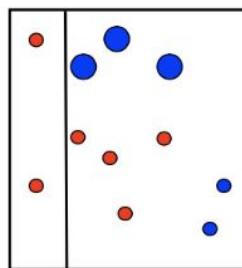
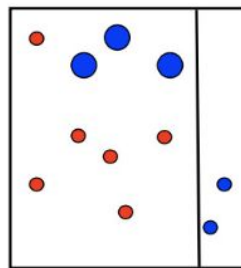
$t = 2$



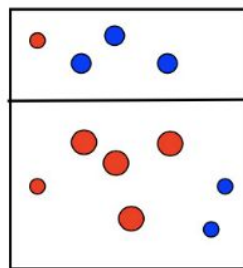
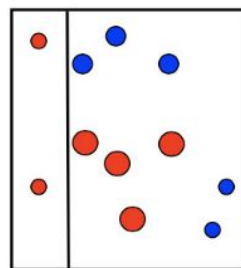
Boosting: intuition



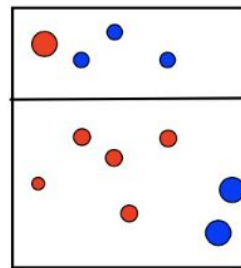
$t = 1$



$t = 2$

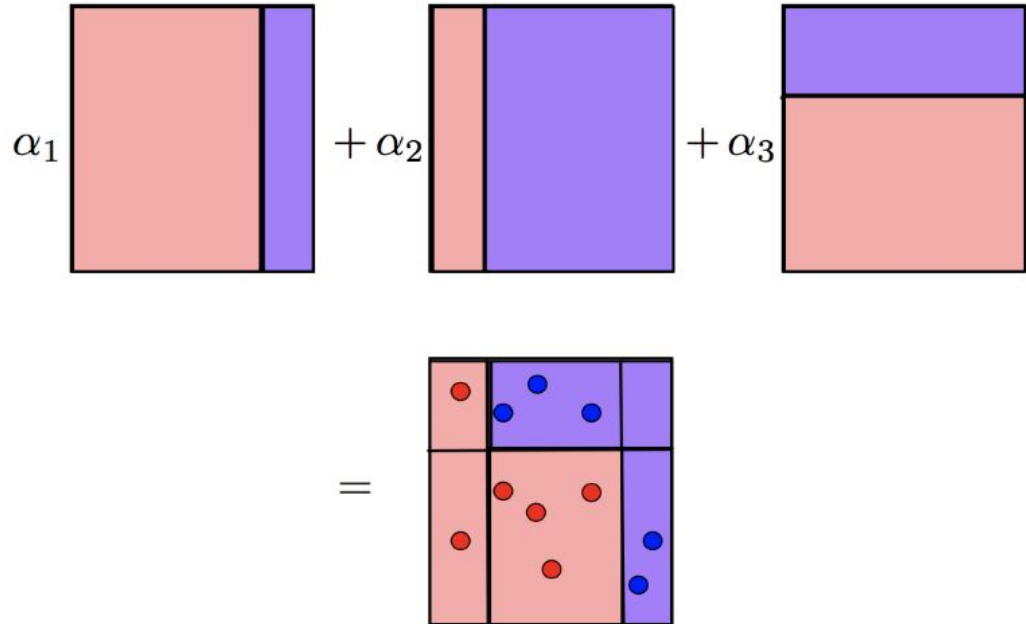
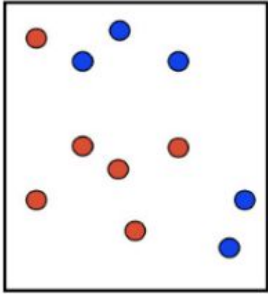


$t = 3$



Boosting: intuition

Binary classification problem.
Models - decision stumps.



Gradient boosting: theory

Denote dataset $\{(x_i, y_i)\}_{i=1, \dots, n}$, loss function $L(y, f)$.

Gradient boosting: theory

Denote dataset $\{(x_i, y_i)\}_{i=1, \dots, n}$, loss function $L(y, f)$.

Optimal model:

$$\hat{f}(x) = \arg \min_{f(x)} L(y, f(x)) = \arg \min_{f(x)} \mathbb{E}_{x,y}[L(y, f(x))]$$

Gradient boosting: theory

Denote dataset $\{(x_i, y_i)\}_{i=1, \dots, n}$, loss function $L(y, f)$.

Optimal model:

$$\hat{f}(x) = \arg \min_{f(x)} L(y, f(x)) = \arg \min_{f(x)} \mathbb{E}_{x,y}[L(y, f(x))]$$

Let it be from parametric family: $\hat{f}(x) = f(x, \hat{\theta})$,

$$\hat{\theta} = \arg \min_{\theta} \mathbb{E}_{x,y}[L(y, f(x, \theta))]$$

Gradient boosting: theory

$$\hat{f}(x) = \sum_{i=0}^{t-1} \hat{f}_i(x),$$

$$(\rho_t, \theta_t) = \arg \min_{\rho, \theta} \mathbb{E}_{x,y} [L(y, \hat{f}(x) + \rho \cdot h(x, \theta))],$$

$$\hat{f}_t(x) = \rho_t \cdot h(x, \theta_t)$$

Gradient boosting: theory

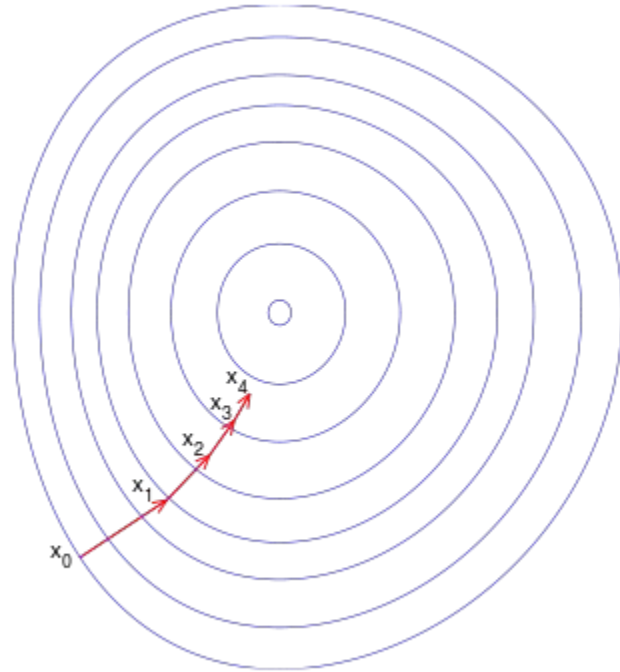
$$\hat{f}(x) = \sum_{i=0}^{t-1} \hat{f}_i(x),$$

$$(\rho_t, \theta_t) = \arg \min_{\rho, \theta} \mathbb{E}_{x,y} [L(y, \hat{f}(x) + \rho \cdot h(x, \theta))],$$

$$\hat{f}_t(x) = \rho_t \cdot h(x, \theta_t)$$

What if we could use gradient descent in *space of our models*?

Gradient boosting: theory



What if we could use gradient descent in *space of our models*?

Gradient boosting: theory

$$\hat{f}(x) = \sum_{i=0}^{t-1} \hat{f}_i(x),$$

$$r_{it} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x)=\hat{f}(x)}, \quad \text{for } i = 1, \dots, n,$$

$$\theta_t = \arg \min_{\theta} \sum_{i=1}^n (r_{it} - h(x_i, \theta))^2,$$

$$\rho_t = \arg \min_{\rho} \sum_{i=1}^n L(y_i, \hat{f}(x_i) + \rho \cdot h(x_i, \theta_t))$$

Gradient boosting: theory

In linear regression case with MSE loss:

$$r_{it} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x)=\hat{f}(x)} = -2(\hat{y}_i - y_i) \propto \hat{y}_i - y_i$$

Gradient boosting: beautiful demo

Great demo:

http://arogozhnikov.github.io/2016/06/24/gradient_boosting_explained.html

Gradient boosting: theory

What we need:

- Data.
- Loss function and its gradient.
- Family of algorithms (with constraints on hyperparameters if necessary).
- Number of iterations M .
- Initial value (GBM by Friedman): constant.

Gradient boosting: example

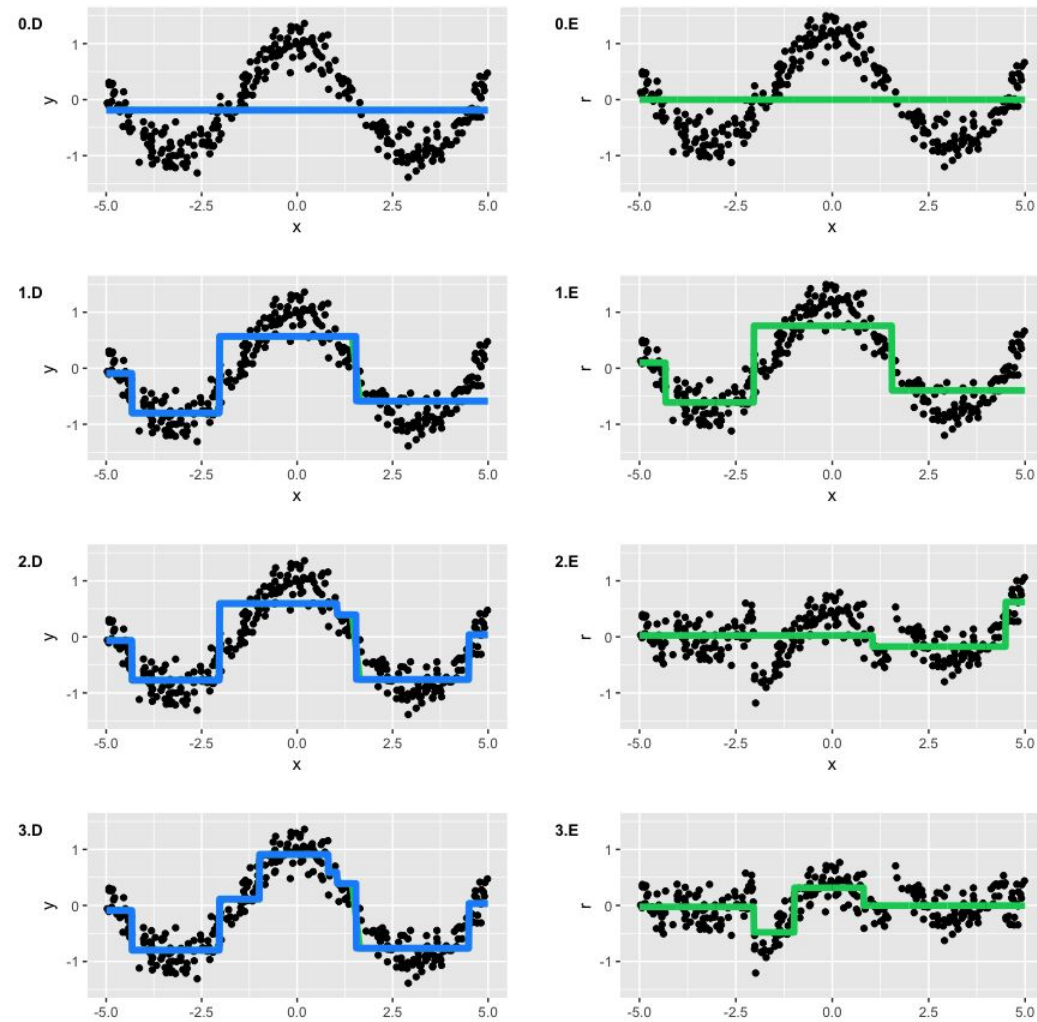
What we need:

- Data: toy dataset $y = \cos(x) + \epsilon, \epsilon \sim \mathcal{N}(0, \frac{1}{5}), x \in [-5, 5]$
- Loss function: MSE
- Family of algorithms: decision trees with depth 2
- Number of iterations $M = 3$
- Initial value: just mean value

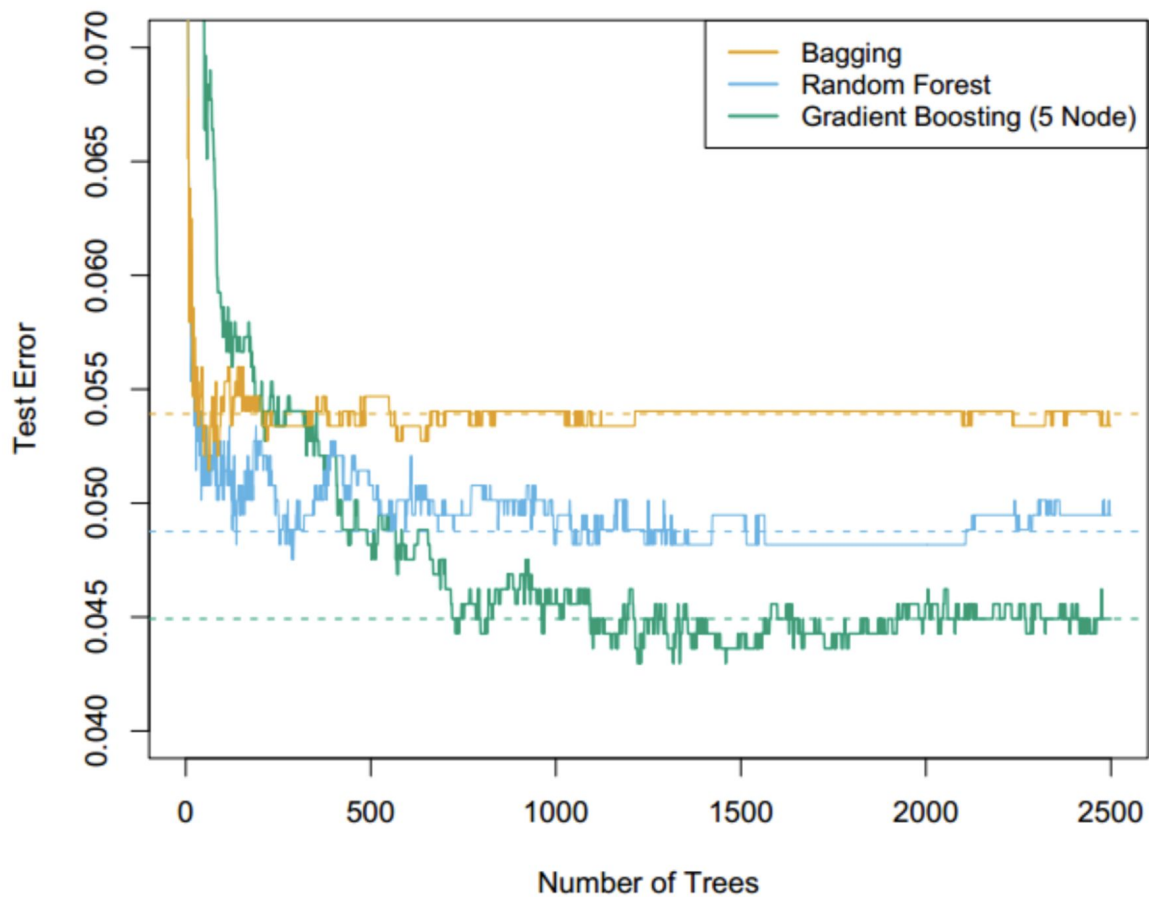
Gradient boosting: example

Left: full ensemble on each step.

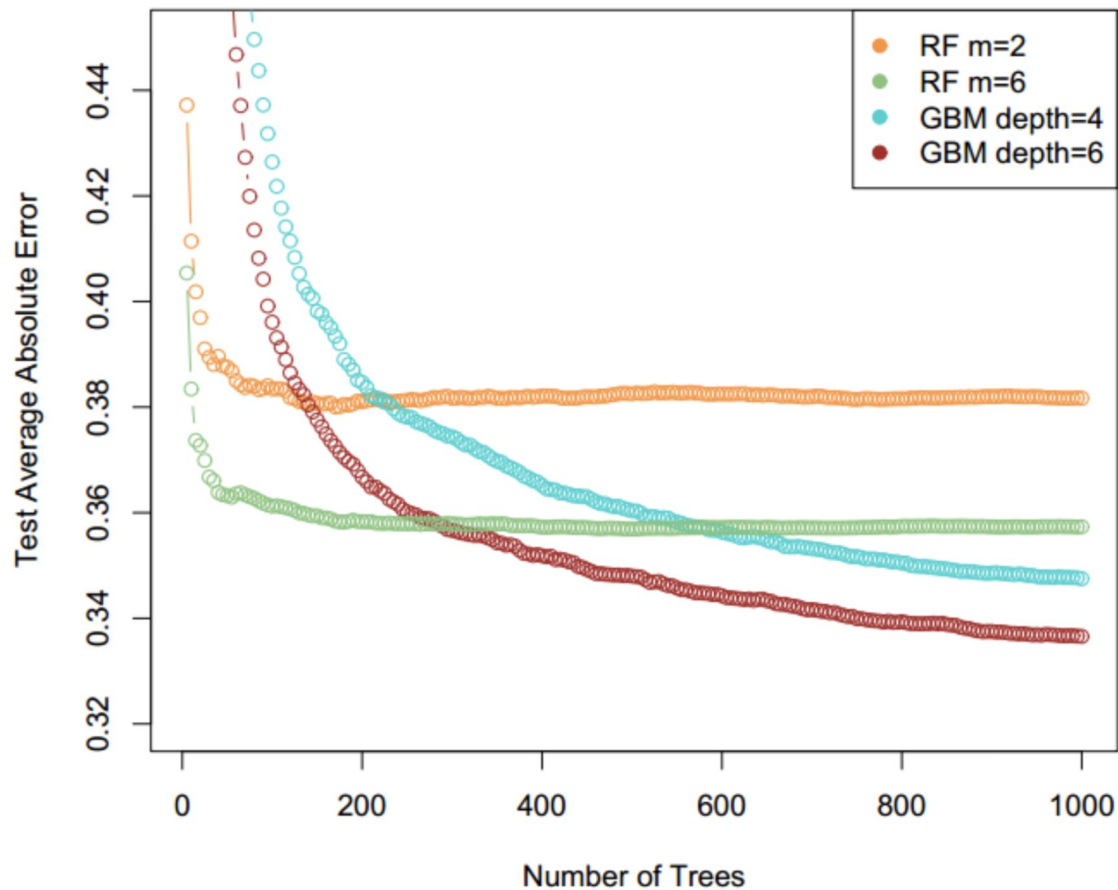
Right: additional tree decisions.



Spam Data

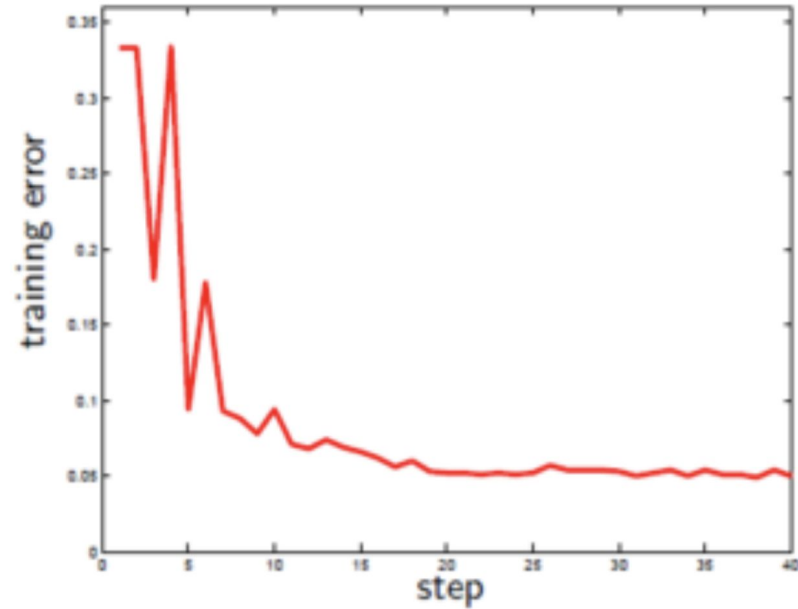
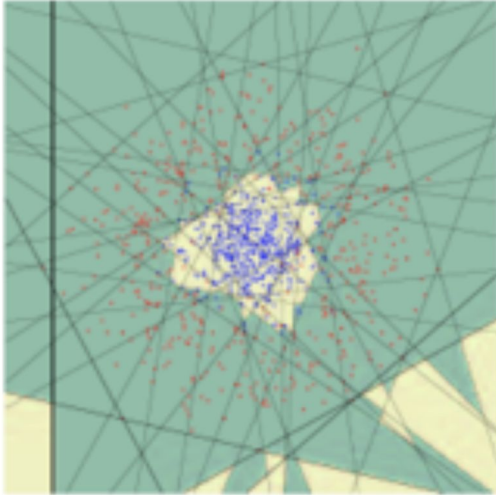


California Housing Data



Boosting with linear classification methods

$t = 40$



Technical side: training in parallel

Which of the ensembling methods could be parallelized?

Technical side: training in parallel

Which of the ensembling methods could be parallelized?

- Random Forest: parallel on the forest level (all trees are independent)

Technical side: training in parallel

Which of the ensembling methods could be parallelized?

- Random Forest: parallel on the forest level (all trees are independent)
- Gradient boosting: parallel on one tree level

Recap: ensembling methods

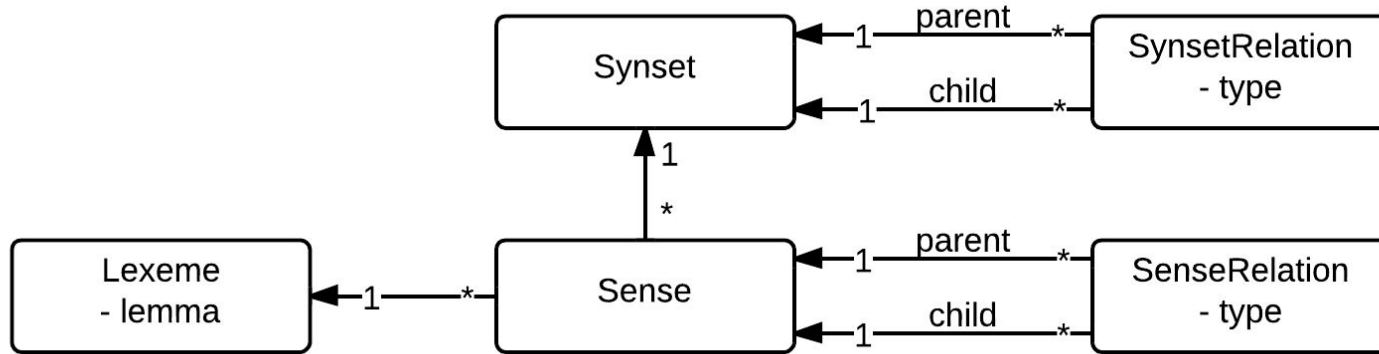
1. Bagging.
2. Random subspace method (RSM).
3. Bagging + RSM + Decision trees = Random Forest.
4. Gradient boosting.
5. Stacking.
6. Blending.

Great demo: http://arogozhnikov.github.io/2016/06/24/gradient_boosting_explained.html

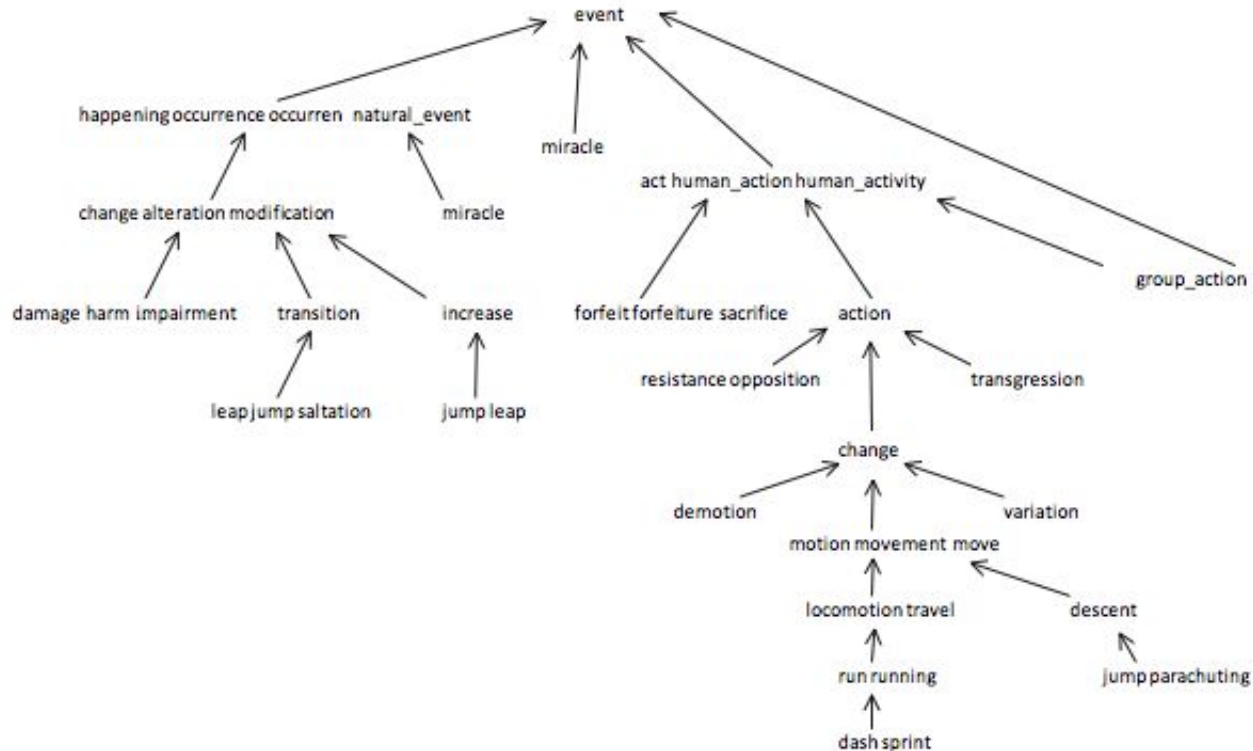
Offtop: words representations

How to represent text in a computer?

Use a taxonomy like WordNet that has hypernyms (is-a) relationships and synonym sets



How to represent text in a computer: WordNet



Discrete representations: problems

- Missing new words
- Subjective
- Requires human labor to create and adapt
- Hard to compute accurate word similarity

Discrete representations: one-hot encoding

"a"	"abbreviations"		"zoology"	"zoom"
1	0		0	0
0	1		0	1
0	0		0	0
⋮	⋮	⋮	⋮	⋮
0	0		0	0
0	0		1	0
0	0		0	1

The diagram illustrates the TF-IDF formula with the following components and annotations:

- Summation:** \sum_w is annotated with "The more query words we match, the better. Σ over the vocabulary".
- Term Frequency (TF):** $tf_{w,Q}$ is circled and annotated with "If word is repeated in the query, it's probably important".
- Document Frequency (IDF) Denominator:** The term $tf_{w,D}$ in the denominator is circled and annotated with "Repetitions of query words in the document \rightarrow good". The entire denominator $tf_{w,D} + \frac{k|D|}{avg|D|}$ is bracketed and annotated with "Repetitions of same word less important than different words. Except in very long documents".
- Inverse Document Frequency (IDF):** The term $\log \frac{|C|}{df_w}$ is circled and annotated with "Rare words more important".

$$s(Q,D) = \sum_w tf_{w,Q} \cdot \frac{tf_{w,D}}{tf_{w,D} + \frac{k|D|}{avg|D|}} \cdot \log \frac{|C|}{df_w}$$

TF - term frequency

IDF - Inversed Document Frequency

TF-IDF: make it simple

$$\text{tf}(\text{"this"}, d_1) = \frac{1}{5} = 0.2$$

$$\text{tf}(\text{"this"}, d_2) = \frac{1}{7} \approx 0.14$$

$$\text{idf}(\text{"this"}, D) = \log\left(\frac{2}{2}\right) = 0$$



$$\text{tfidf}(\text{"this"}, d_1, D) = 0.2 \times 0 = 0$$

$$\text{tfidf}(\text{"this"}, d_2, D) = 0.14 \times 0 = 0$$



Word 'this' is not very
informative

Document 1		Document 2	
Term	Term Count	Term	Term Count
this	1	this	1
is	1	is	1
a	2	another	2
sample	1	example	3

One of the most successful ideas of statistical NLP:

“You shall know a word by the company it keeps”

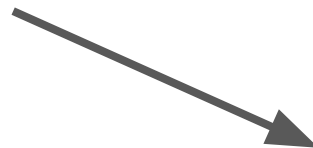
(J. R. Firth 1957: 11)

Words cooccurrences

Finding N-grams in a text



Word-document
cooccurrence matrix

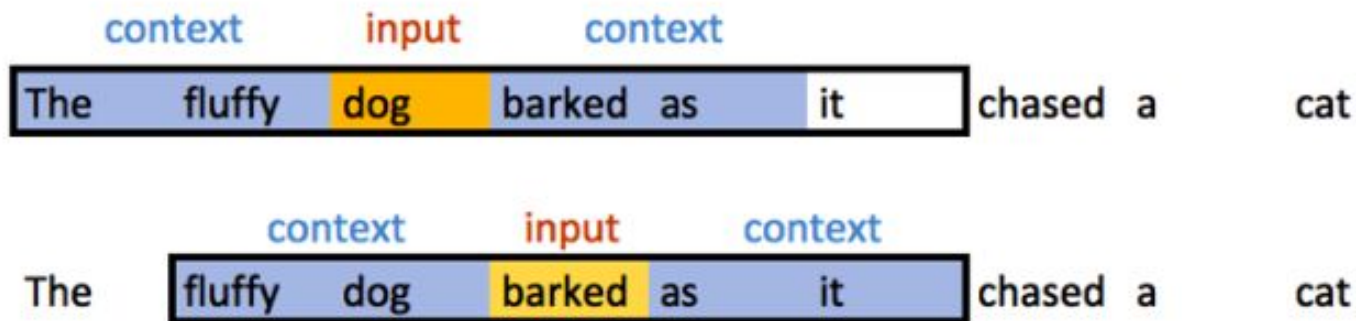


Window around
each word

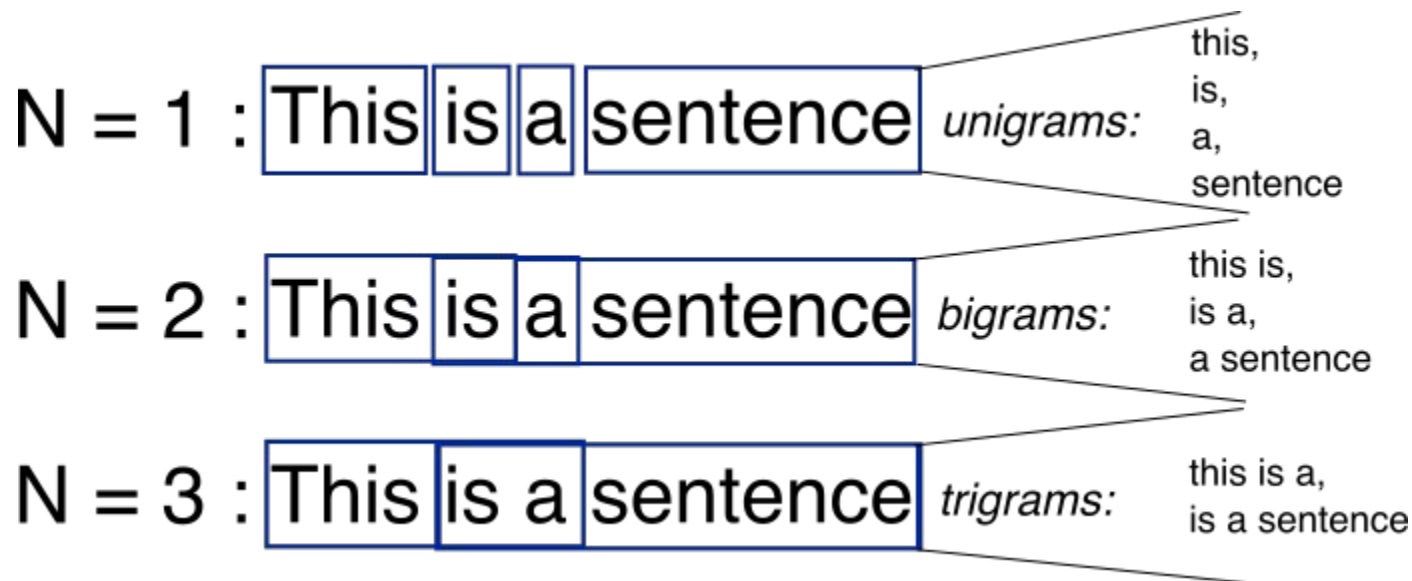
Word-document cooccurrence matrix

$$X = \begin{matrix} & \begin{matrix} I & like & enjoy & deep & learning & NLP & flying & . \end{matrix} \\ \begin{matrix} I \\ like \\ enjoy \\ deep \\ learning \\ NLP \\ flying \\ . \end{matrix} & \left[\begin{array}{cccccccc} 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{array} \right] \end{matrix}$$

Words cooccurrences: sliding window



Words cooccurrences: n-grams



Cooccurrence vectors: problems

- Increase in size with vocabulary
- Very high dimensional: require a lot of storage
- Subsequent classification models have sparsity issues



Models are less robust

More interesting approaches
coming in next classes.
Stay tuned!

