

INTRO TO DATA SCIENCE

COURSE REVIEW & WHERE TO GO NEXT

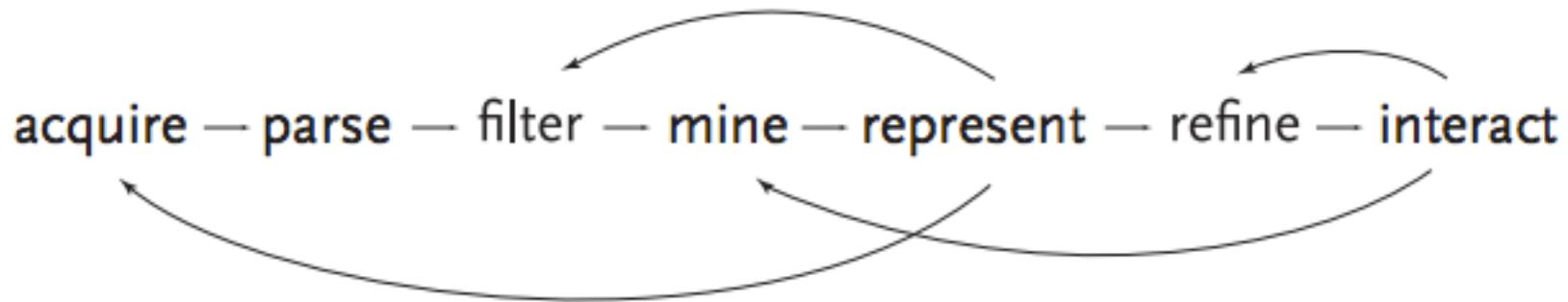
I. WHERE HAVE WE BEEN?

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

- *Data acquisition and preparation*
- *Exploratory data analysis*
- *Supervised learning:*
 - *kNN*
 - *Linear, multiple, & polynomial regression*
 - *Decision trees & random forests*
 - *logistic regression*
 - *Naïve Bayes*
 - *Support Vector Machines (SVM)*

- *Unsupervised learning:*
 - *K-means clustering*
 - *PCA/SVD for dimensionality reduction*
- *Model evaluation, confusion matrix, ROC curves, AUC*
- *APIs & web scraping*
- *Recommender systems*
- *Databases and SQL*
- *Text mining and natural language processing*

THE DATA SCIENCE WORKFLOW



NOTE

This diagram illustrates
the *iterative* nature of
problem solving

GIT COMMANDS

Main

git clone – clone a repo

git status – get status

git add – add changes to be pushed

git commit – commit the change with a comment

git push – push the change to github

git pull – pull remote changes from github

Others

git branch – see all branches

git checkout – checkout a branch

git merge – merge in another branch

git stash – stash changes

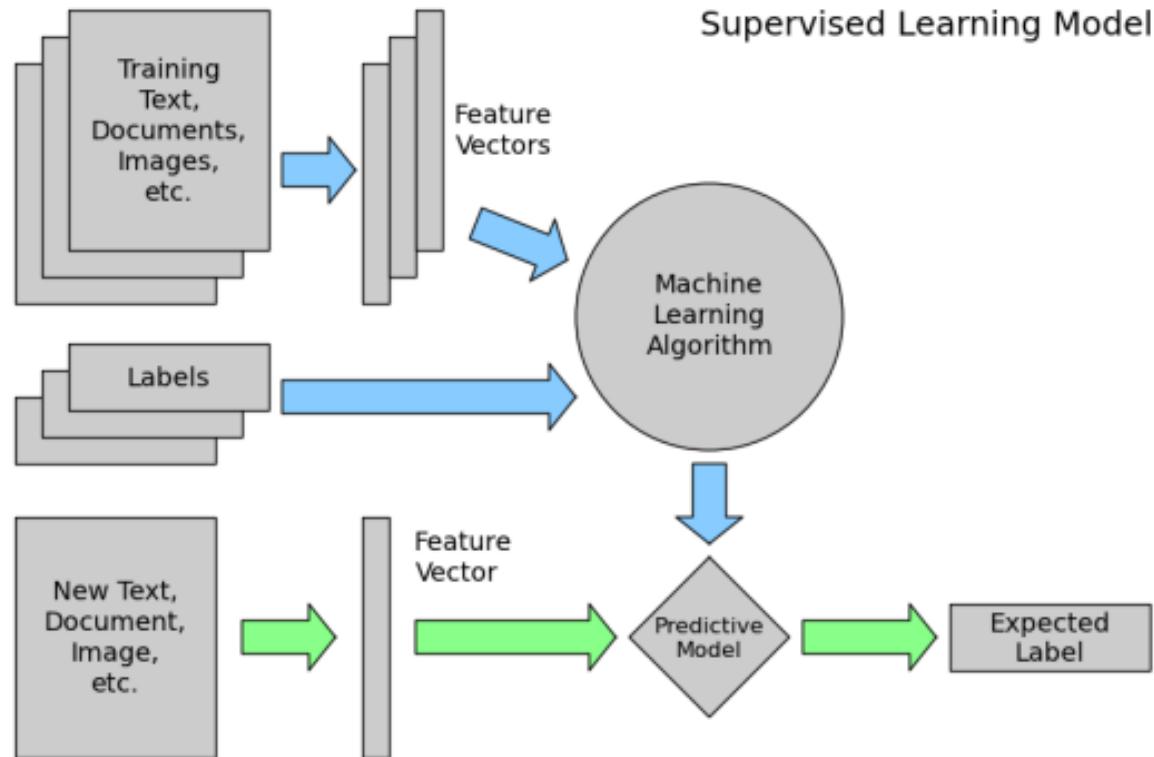
pull request – remote changes requested to be merged in



PEP 20: THE ZEN OF PYTHON

Beautiful is better than ugly.
Explicit is better than implicit.
Simple is better than complex.
Complex is better than complicated.
Flat is better than nested.
Sparse is better than dense.
Readability counts.
Special cases aren't special enough to break the rules.
Although practicality beats purity.
Errors should never pass silently.
Unless explicitly silenced.
In the face of ambiguity, refuse the temptation to guess.
There should be one-- and preferably only one --obvious way to do it.
Although that way may not be obvious at first unless you're Dutch.
Now is better than never.
Although never is often better than *right* now.
If the implementation is hard to explain, it's a bad idea.
If the implementation is easy to explain, it may be a good idea.

How does supervised learning “work”?



Suppose we want to predict the color of the gray dot.

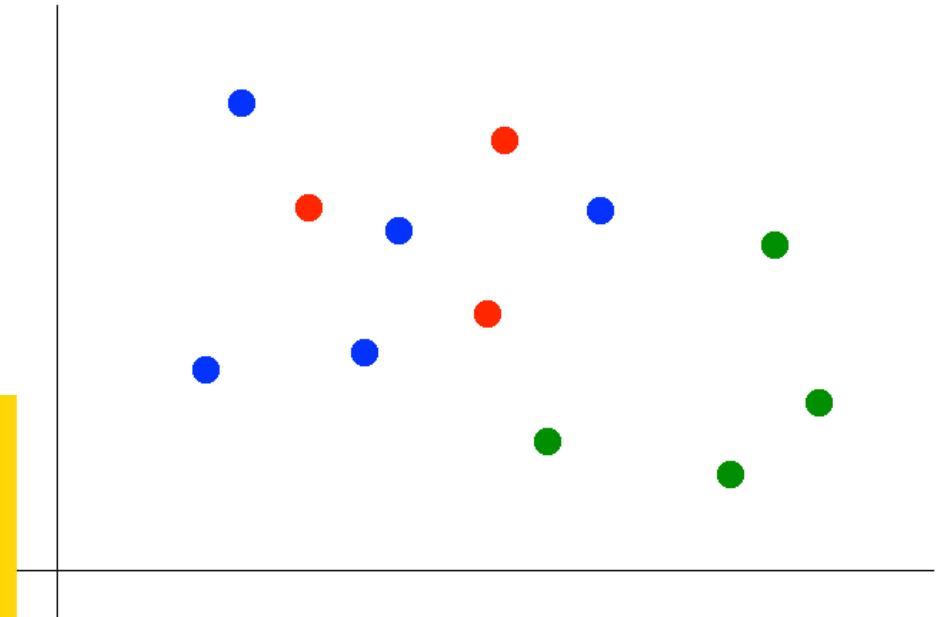
1) *Pick a value for k.*

2) *Find colors of k nearest neighbors.*

3) *Assign the most common color to the gray dot.*

NOTE:

Our definition of “nearest” implicitly uses the *Euclidean distance function*.



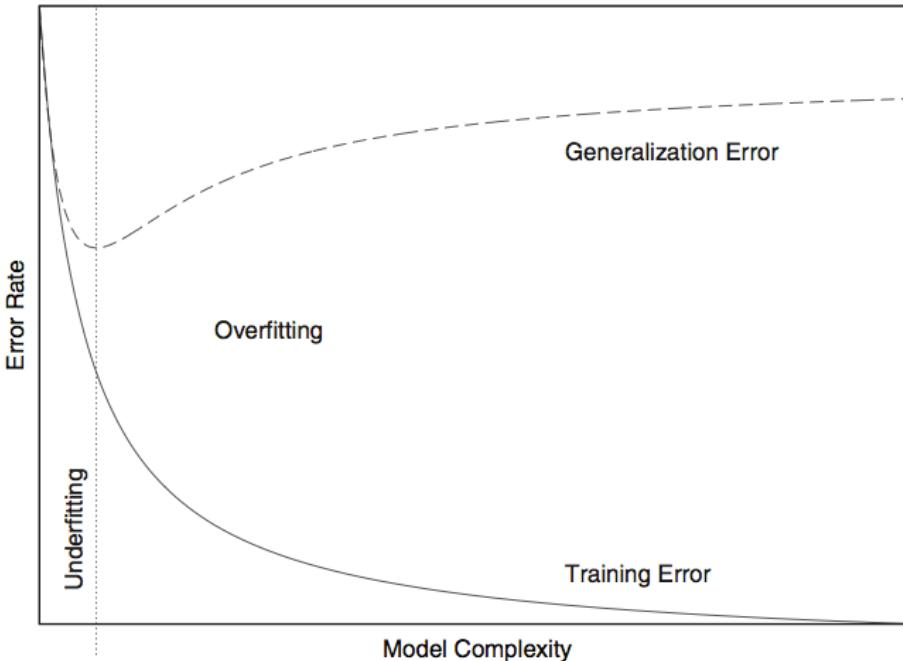
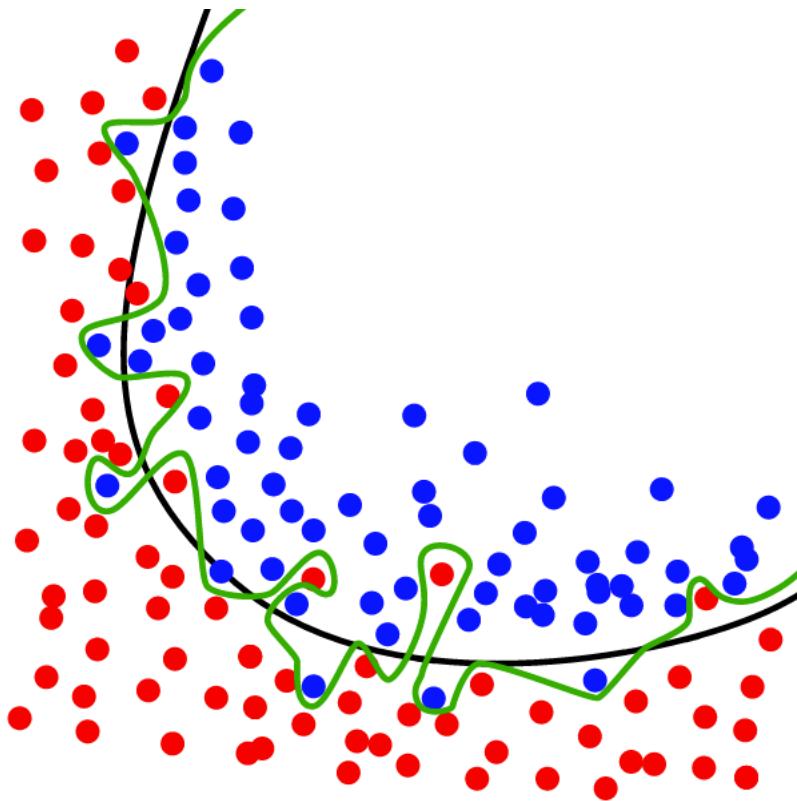


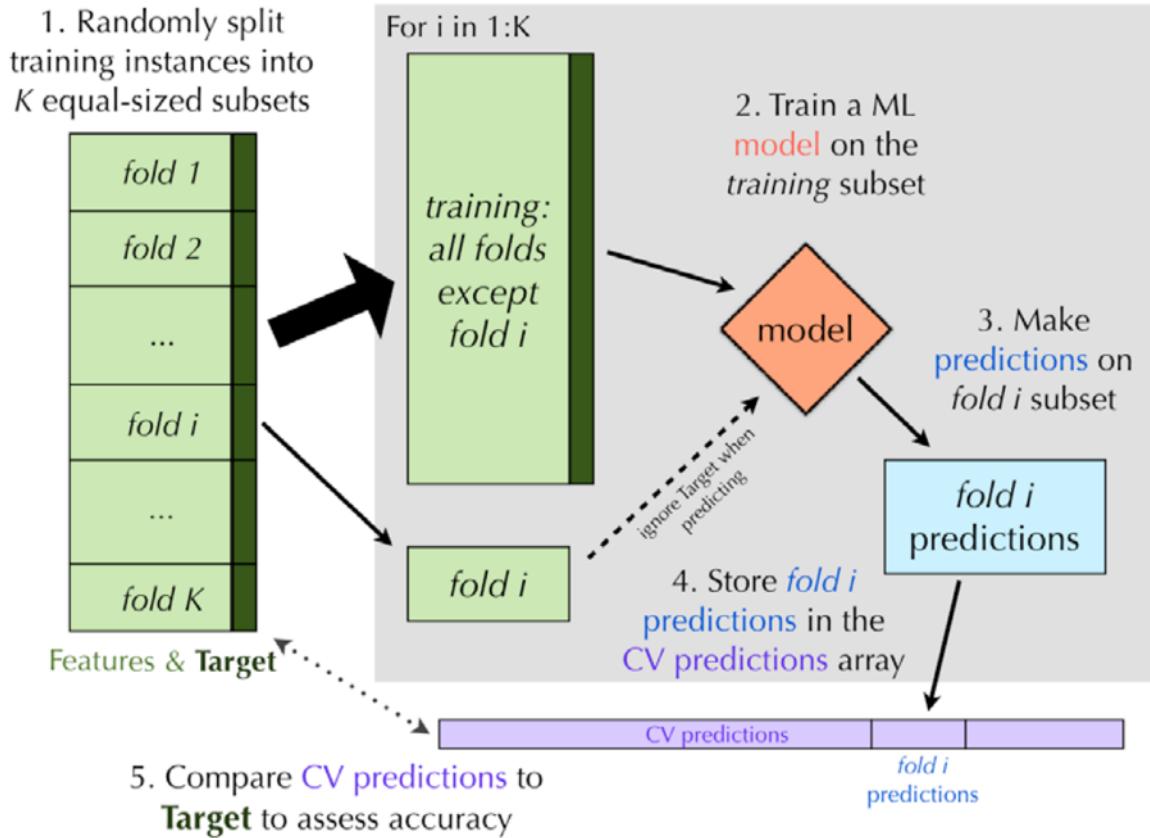
FIGURE 18-1. Overfitting: as a model becomes more complex, it becomes increasingly able to represent the training data. However, such a model is overfitted and will not generalize well to data that was not used during training.

source: *Data Analysis with Open Source Tools*, by Philipp K. Janert. O'Reilly Media, 2011.



Dataset	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	<u>Accuracy</u>
1	Test	Train	Train	Train	Train	$k_1 \%$
2	Train	Test	Train	Train	Train	$k_2 \%$
3	Train	Train	Test	Train	Train	$k_3 \%$
4	Train	Train	Train	Test	Train	$k_4 \%$
5	Train	Train	Train	Train	Test	$k_5 \%$

$$5\text{-Fold Generalization Error} = (k_1 + k_2 + k_3 + k_4 + k_5) / 5$$



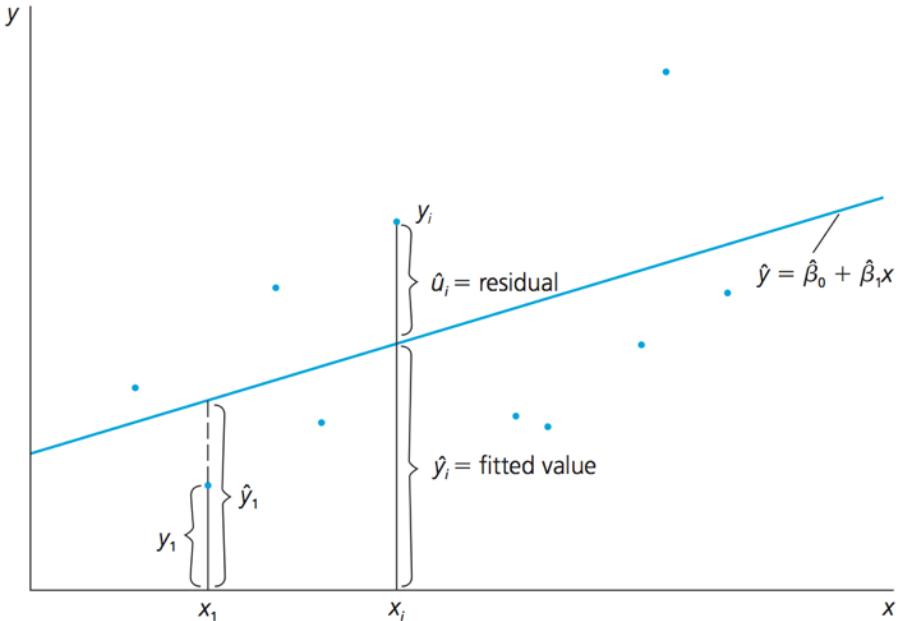
Q: How do we fit a regression model to a dataset?

A: In theory, minimize the sum of the squared residuals (OLS).

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

$$\hat{u}_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i.$$

$$\sum_{i=1}^n \hat{u}_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2,$$



This tradeoff is regulated by a hyperparameter λ , which we've already seen:

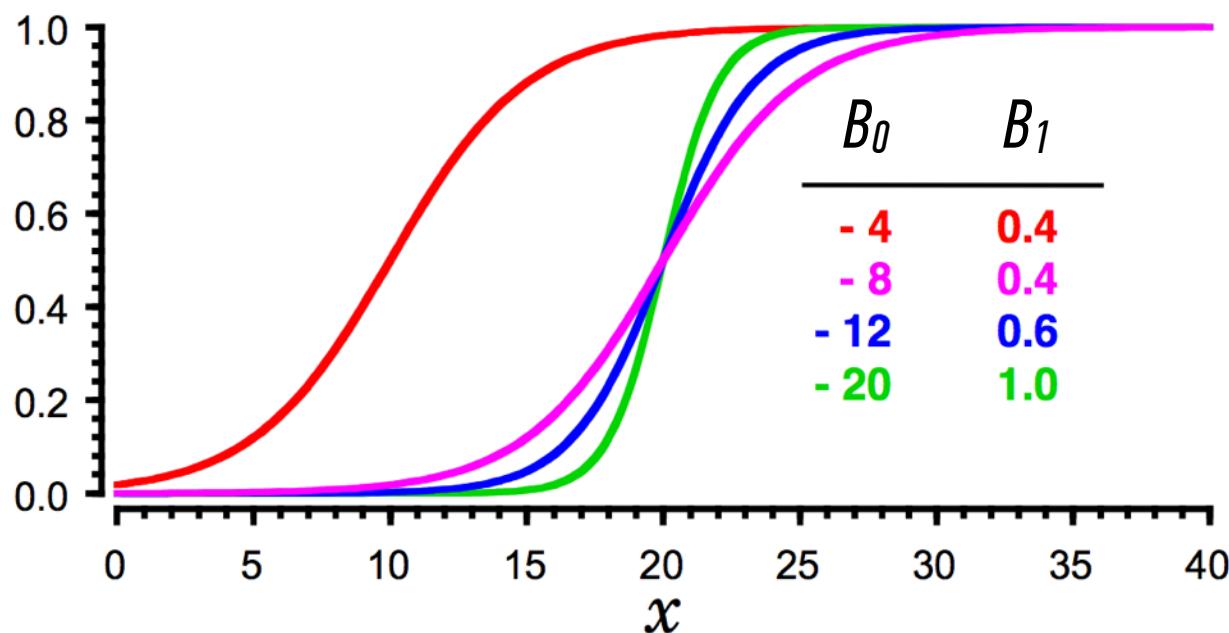
L1 regularization: $y = \sum \beta_i x_i + \varepsilon$ st. $\sum |\beta_i| < \lambda$

L2 regularization: $y = \sum \beta_i x_i + \varepsilon$ st. $\sum \beta_i^2 < \lambda$

So regularization represents a method to trade away some variance for a little bias in our model, thus achieving a better overall fit.

LOGISTIC REGRESSION

$$F(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$



When $\beta_0 + \beta_1 x = 0$, then $F(x) = 0.5$, which is the inflection point on all these curves.

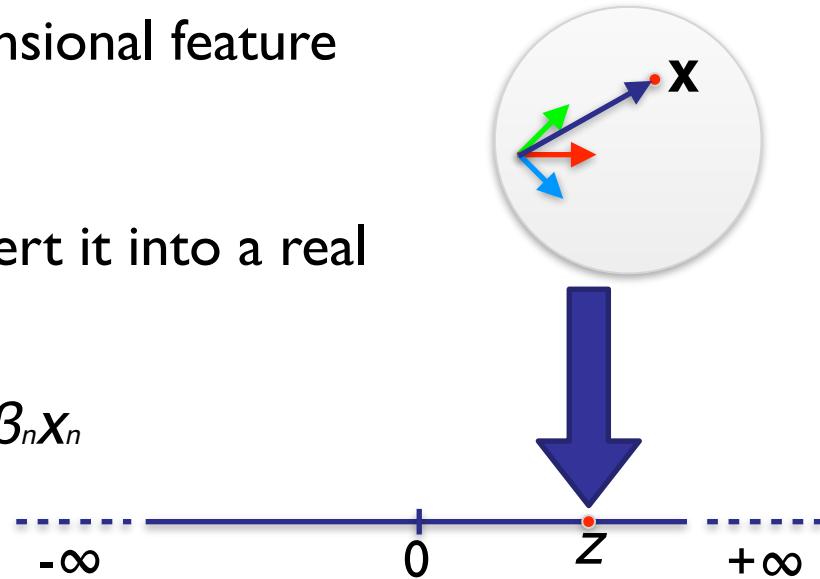
LOGISTIC REGRESSION

- I. Model consists of a vector β in n-dimensional feature space

$$\beta = \beta_1 + \beta_2 + \dots + \beta_n$$

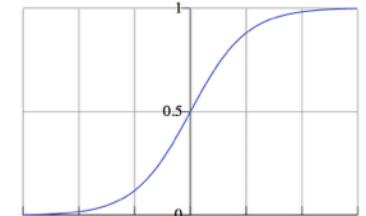
2. For a point x , project it onto β to convert it into a real number z in the range $-\infty$ to $+\infty$

$$z = \alpha + \beta \cdot x = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$



3. Map z to the range 0 to 1 using the logistic function

$$p = 1 / (1 + e^{-z})$$



Notice if we take the logarithm of the odds, we return a linear equation

$$\log\left(\frac{\pi}{1-\pi}\right) = \log(e^{\beta_0 + \beta_1 x}) = \beta_0 + \beta_1 x$$

This simple relationship between the odds ratio and the parameter β is what makes logistic regression such a powerful tool.

*This term is the **posterior probability** of C. It represents the probability of a record belonging to class C after the data is taken into account.*

$$P(\text{class } C | \{x_i\}) = \frac{P(\{x_i\} | \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

The goal of any Bayesian computation is to find (“learn”) the posterior distribution of a particular variable.

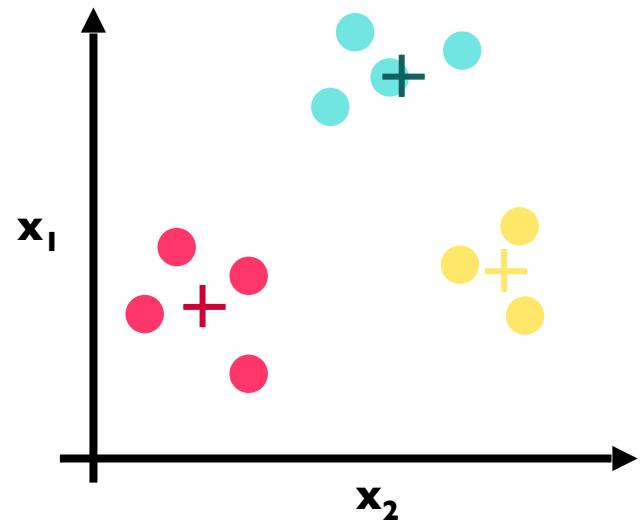
1) choose k initial centroids (note that k is an input)

2) for each point:

- find distance to each centroid
- assign point to nearest centroid

3) recalculate centroid positions

4) repeat steps 2-3 until stopping criteria met



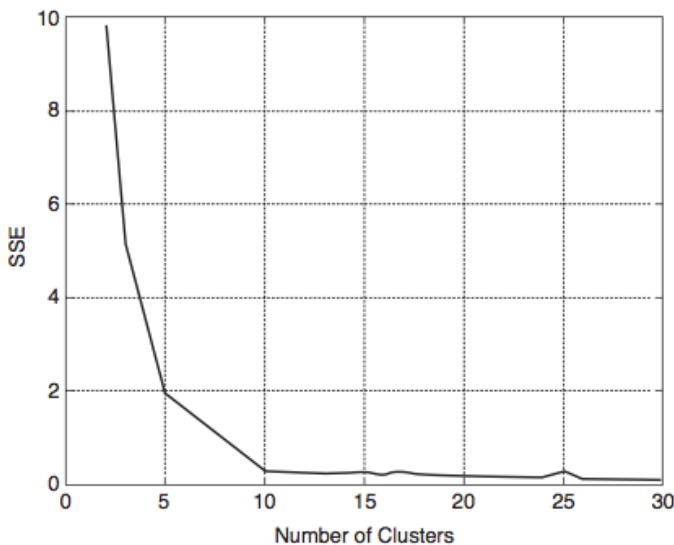


Figure 8.32. SSE versus number of clusters for the data of Figure 8.29.

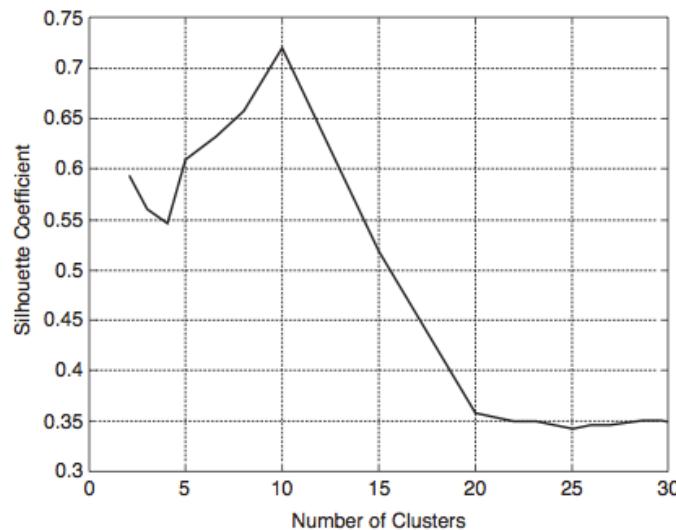
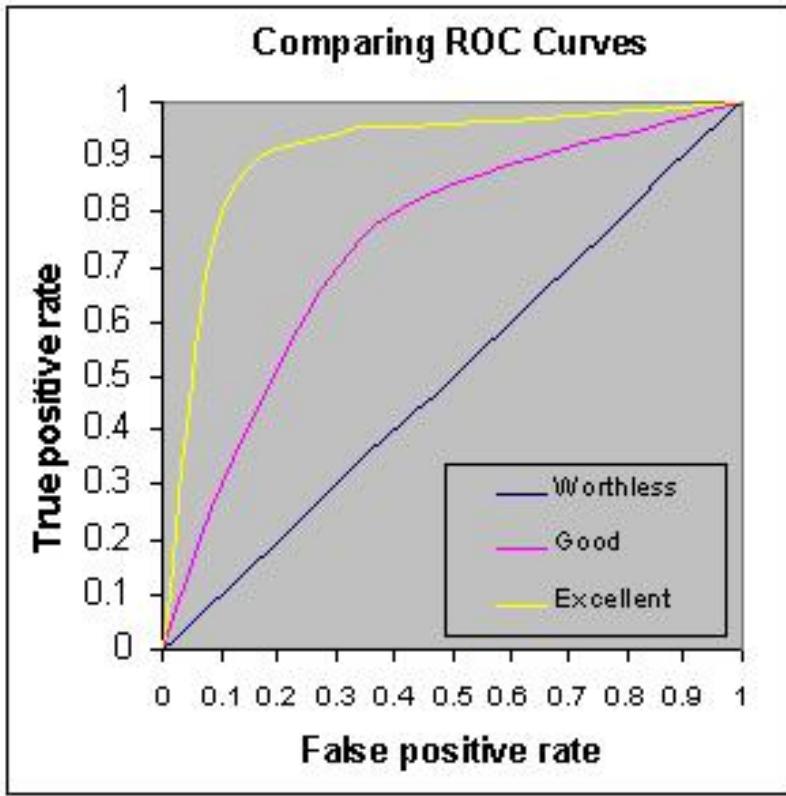
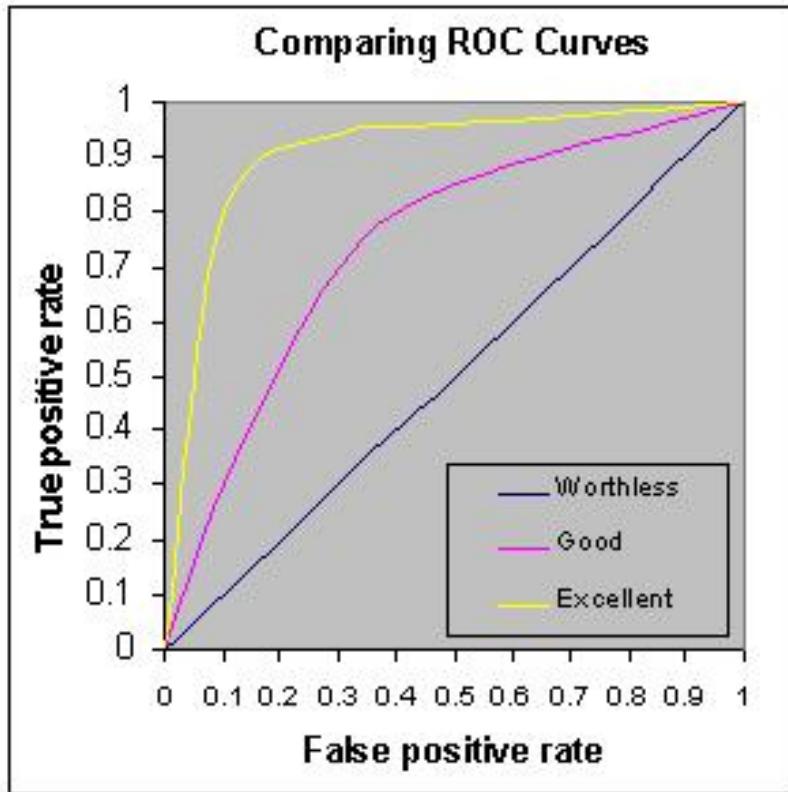


Figure 8.33. Average silhouette coefficient versus number of clusters for the data of Figure 8.29.



ROC Curves show the relationship between the TP Rate and the FP Rate as we vary the decision threshold for the classifier

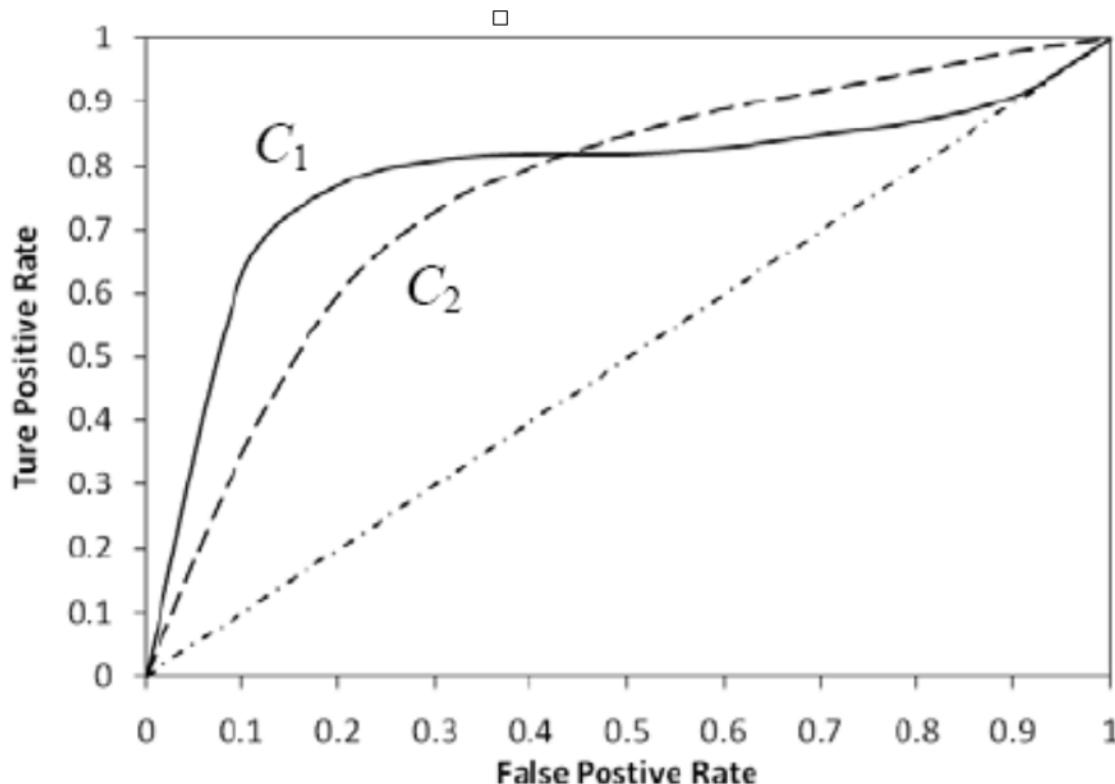
Cut off	TPR (y)	FPR (x)	Cut off	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0



Area Under the Curve (AUC)

We evaluate a classifier by measuring the Area Under the Curve for its ROC curve. The Greater area under the curve, the more effective the classifier.

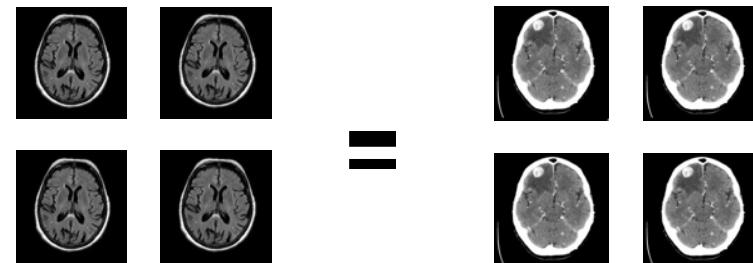
Then for our chosen classifier, we pick an appropriate decision threshold. In general, we pick the decision threshold that gets us closest to the upper left corner

MODEL EVALUATION**Pop Quiz!**

Which of the two classifiers shown (on the same data), C_1 or C_2 , is better and why?

Imbalanced classes can be re-balanced in several ways

1. **Undersampling** the dominant class - remove some the majority class so it has less weight
2. **Oversampling** the minority class - add more of the minority class so it has more weight.
3. **Hybrid** - doing both



*We want our objective function to measure the **gain** in purity from a particular split.*

Table 4.1. The vertebrate data set.

Name	Body Temperature	Skin Cover	Gives Birth	Aquatic Creature	Aerial Creature	Has Legs	Hibernates	Class Label
human	warm-blooded	hair	yes	no	no	yes	no	mammal
python	cold-blooded	scales	no	no	no	no	yes	reptile
salmon	cold-blooded	scales	no	yes	no	no	no	fish
whale	warm-blooded	hair	yes	yes	no	no	no	mammal
frog	cold-blooded	none	no	semi	no	yes	yes	amphibian
komodo dragon	cold-blooded	scales	no	no	no	yes	no	reptile
bat	warm-blooded	hair	yes	no	yes	yes	yes	mammal
pigeon	warm-blooded	feathers	no	no	yes	yes	no	bird
cat	warm-blooded	fur	yes	no	no	yes	no	mammal
leopard	cold-blooded	scales	yes	yes	no	no	no	fish
shark	cold-blooded	scales	no	semi	no	yes	no	reptile
turtle	cold-blooded	scales	no	semi	no	yes	no	bird
penguin	warm-blooded	feathers	no	semi	no	yes	yes	mammal
porcupine	warm-blooded	quills	yes	no	no	yes	yes	fish
eel	cold-blooded	scales	no	yes	no	no	no	mammal
salamander	cold-blooded	none	no	semi	no	yes	yes	amphibian

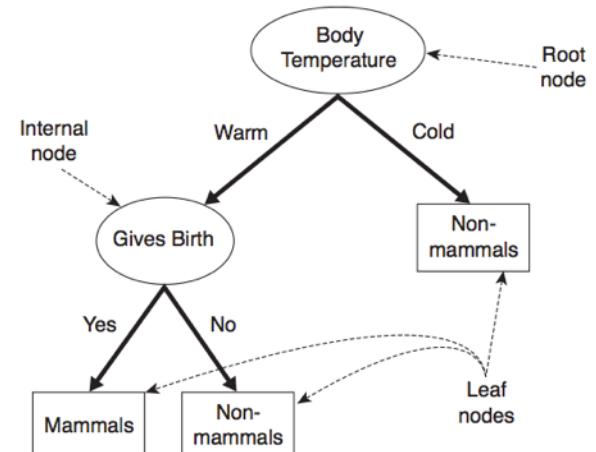
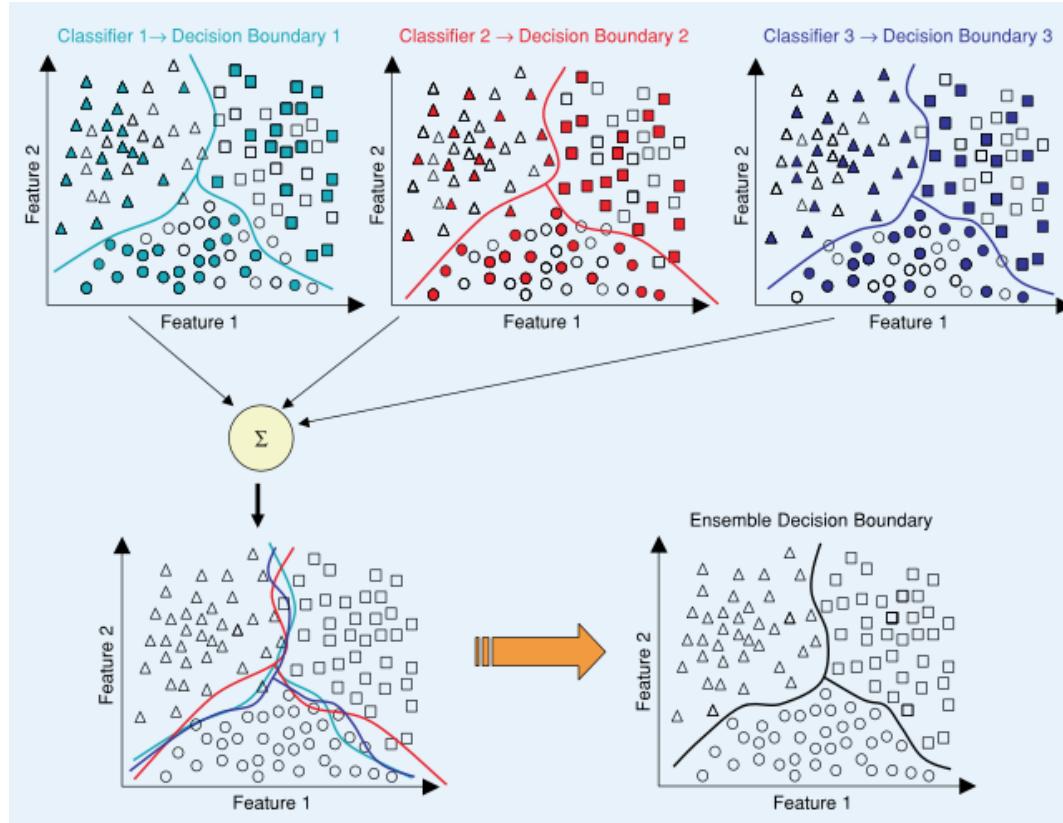


Figure 4.4. A decision tree for the mammal classification problem.

ENSEMBLE METHODS

28



CREATING AN ENSEMBLE PREDICTION

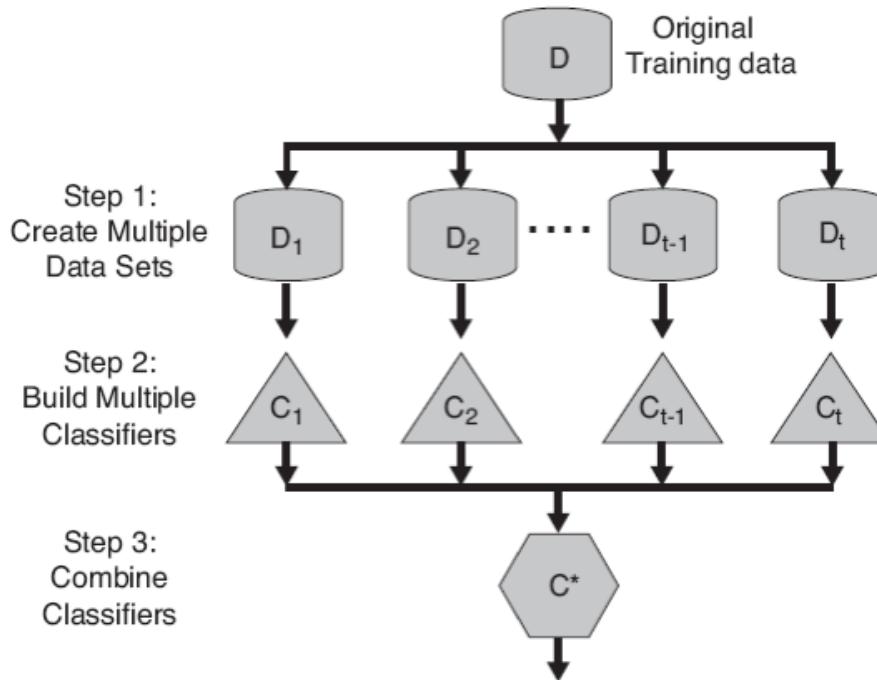
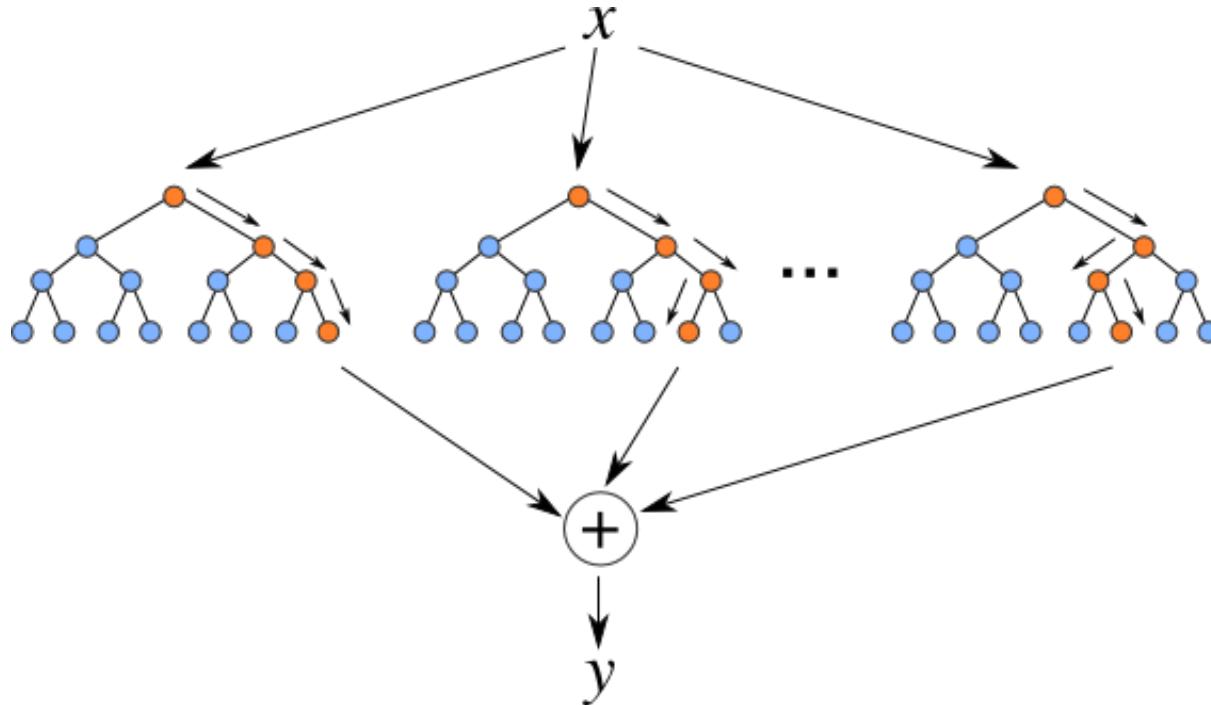
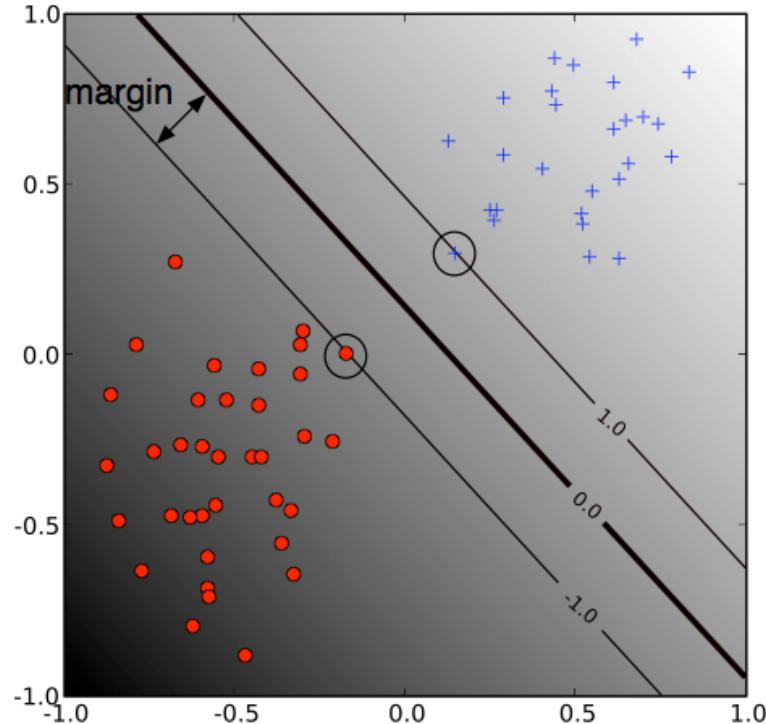


Figure 5.31. A logical view of the ensemble learning method.

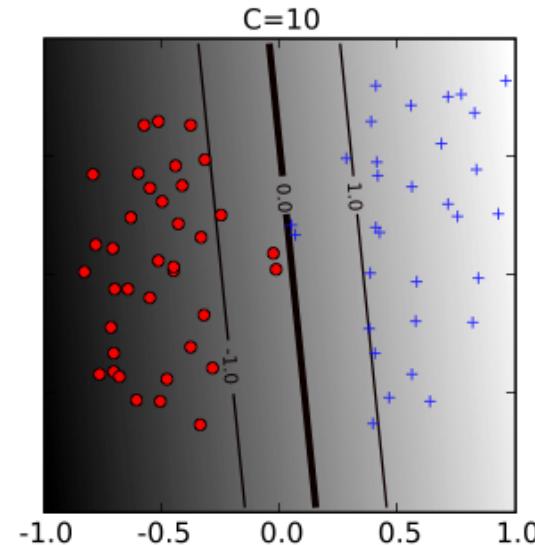
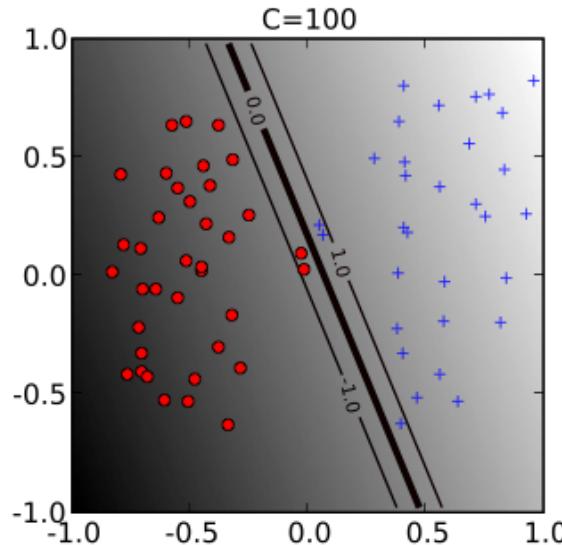


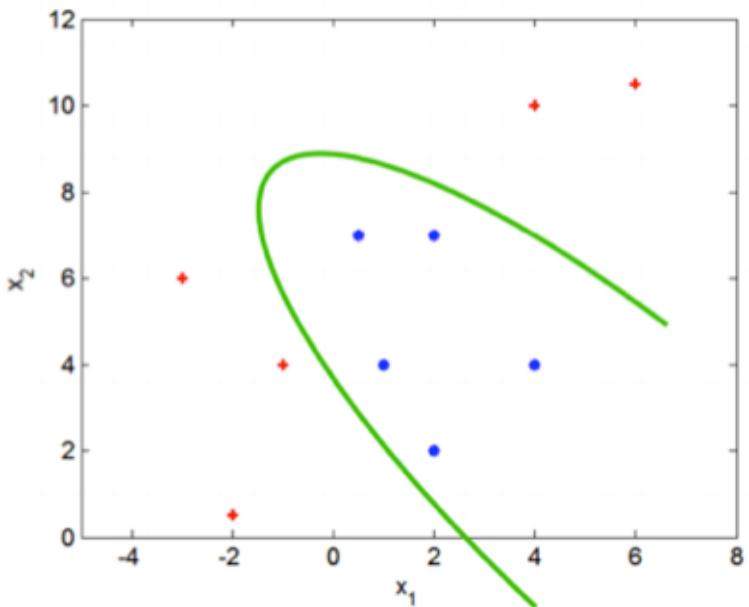
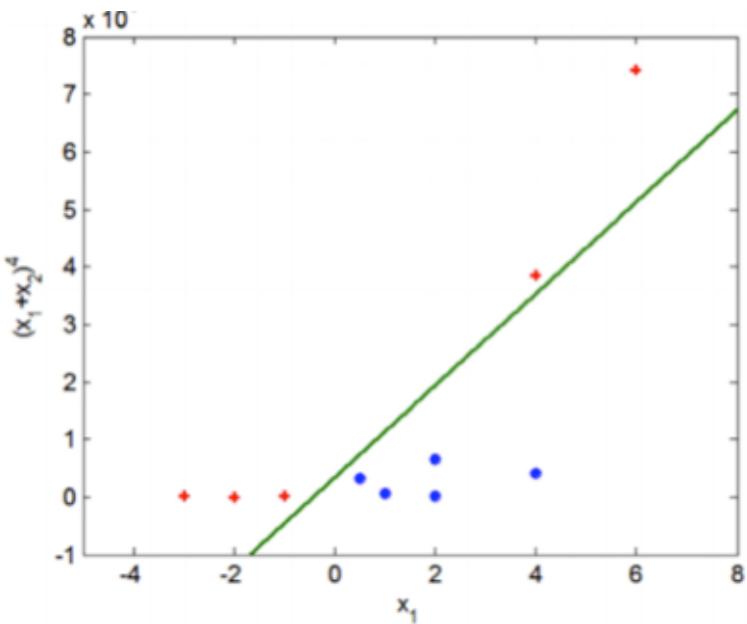
Notice that the margin depends only on the points that are nearest to the decision boundary.

These points are called the support vectors.

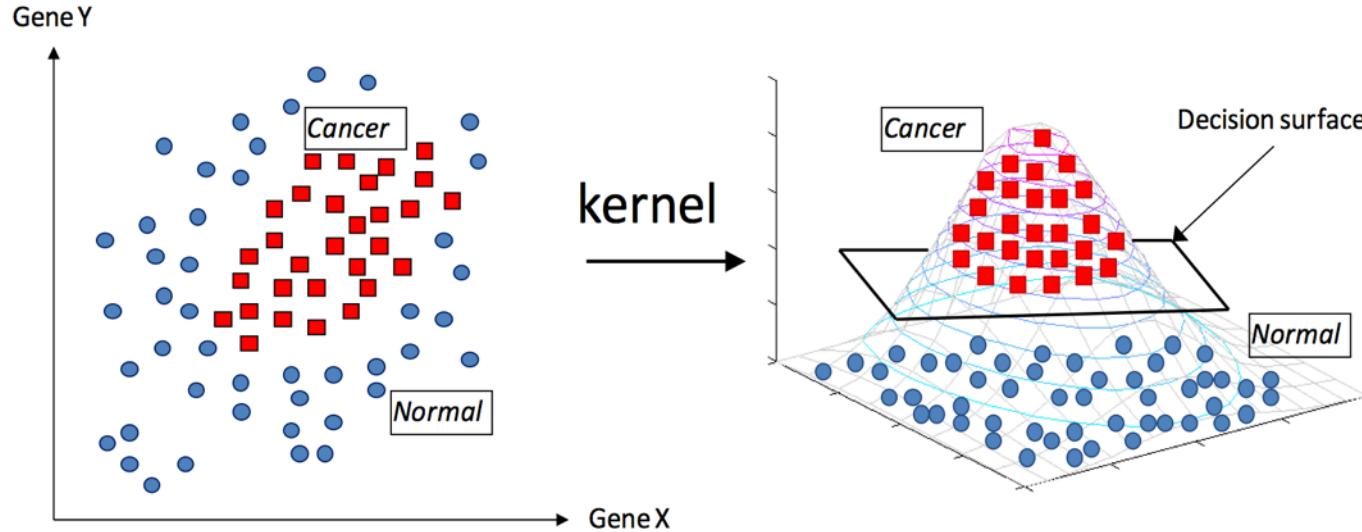


We can trade some training error in order to get a larger margin (remember the bias-variance trade-off)

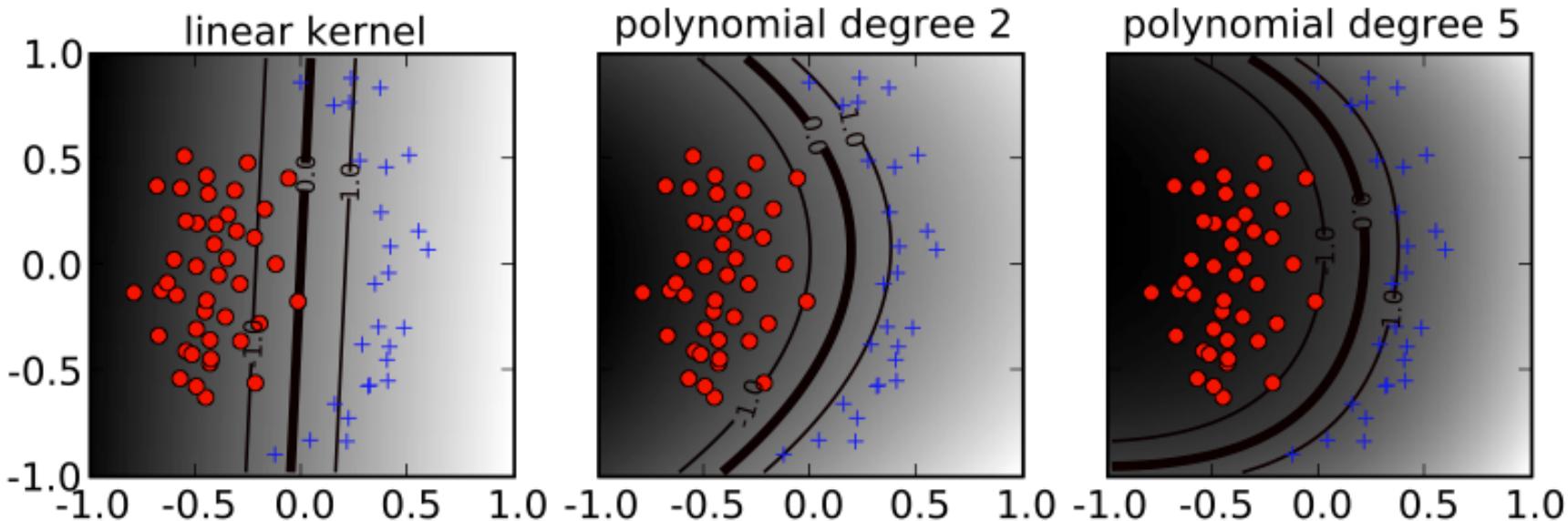


SUPPORT VECTOR MACHINESoriginal feature space K higher-dim feature space K'

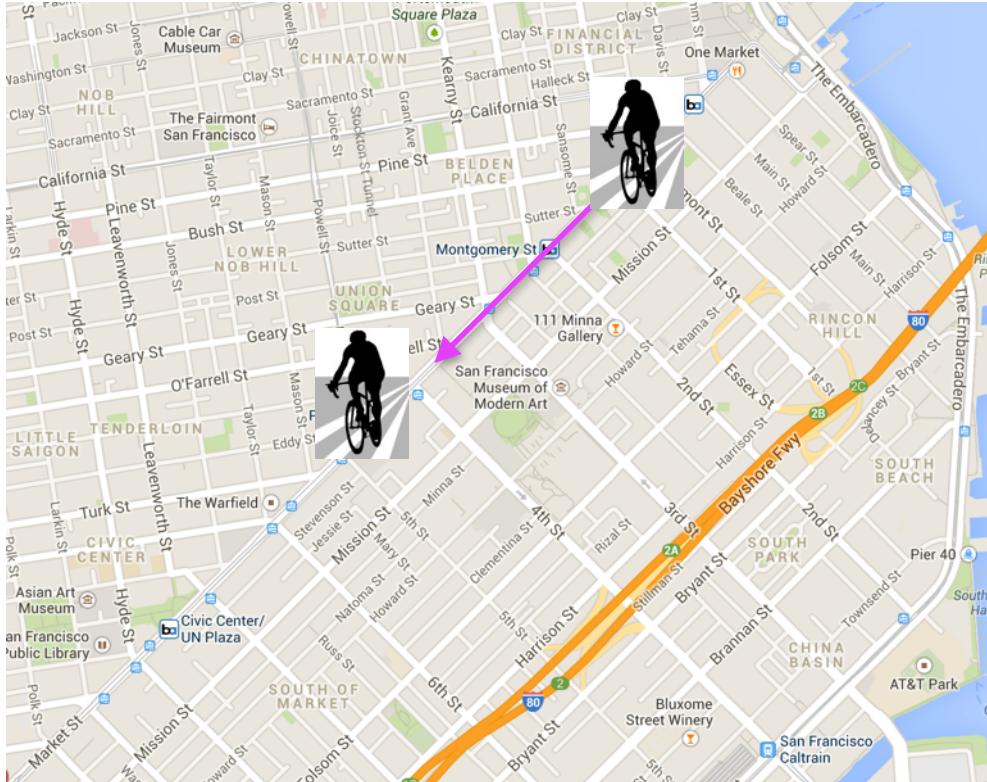
SUPPORT VECTOR MACHINES: THE KERNEL TRICK



If a linear decision boundary cannot be found in the original space, we can map into a higher dimensional space and find the separating surface.

SUPPORT VECTOR MACHINES: NONLINEAR CLASSIFICATION – POLYNOMIAL KERNEL

DIMENSIONALITY REDUCTION



What if we just used
distance down Market St.?

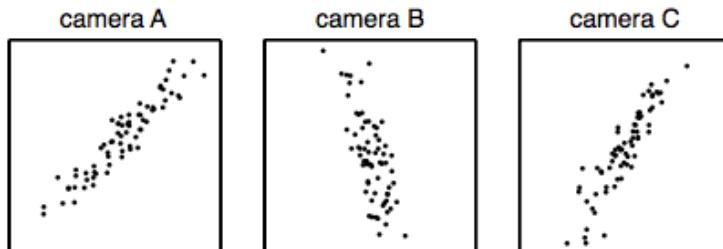
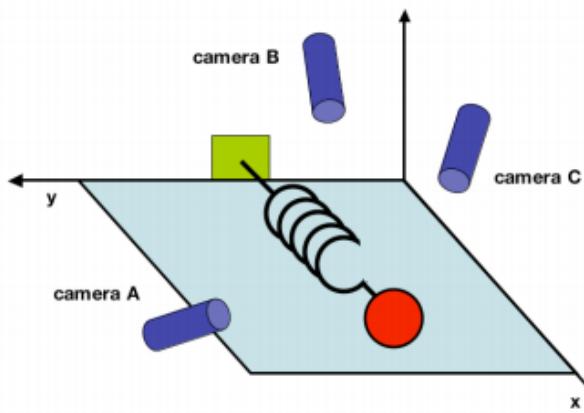
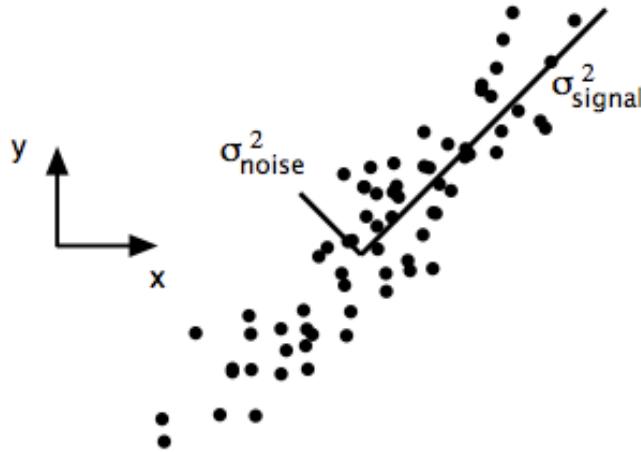


FIG. 1 A toy example. The position of a ball attached to an oscillating spring is recorded using three cameras A, B and C. The position of the ball tracked by each camera is depicted in each panel below.

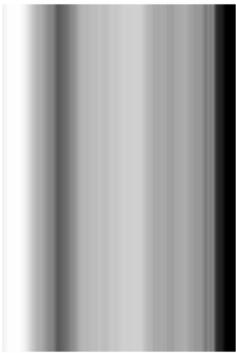


$$SNR = \frac{\sigma_{signal}^2}{\sigma_{noise}^2}.$$

FIG. 2 Simulated data of (x, y) for camera A. The signal and noise variances σ_{signal}^2 and σ_{noise}^2 are graphically represented by the two lines subtending the cloud of data. Note that the largest direction of variance does not lie along the basis of the recording (x_A, y_A) but rather along the best-fit line.

DIMENSIONALITY REDUCTION: PCA

PCs # 0



PCs # 10



PCs # 20



PCs # 30



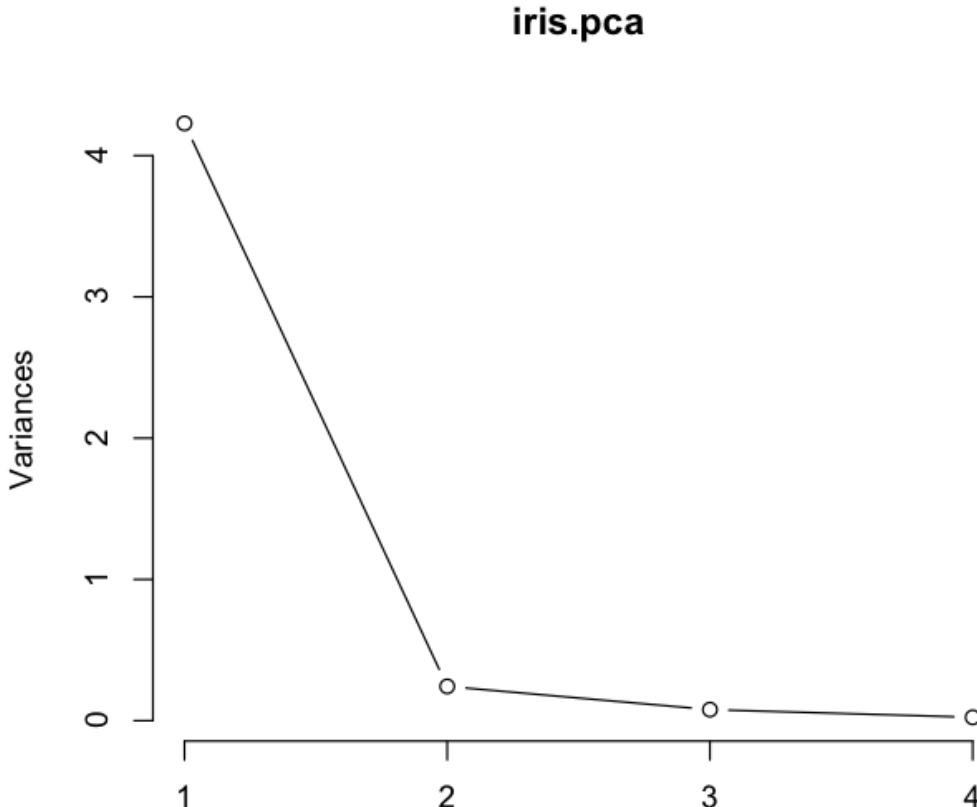
PCs # 40



PCs # 50



PRINCIPAL COMPONENT ANALYSIS



NOTE

Looking at this plot also gives you an idea of how many principal components to keep.

Apply the *elbow test*: keep only those pc's that appear to the left of the elbow in the graph.

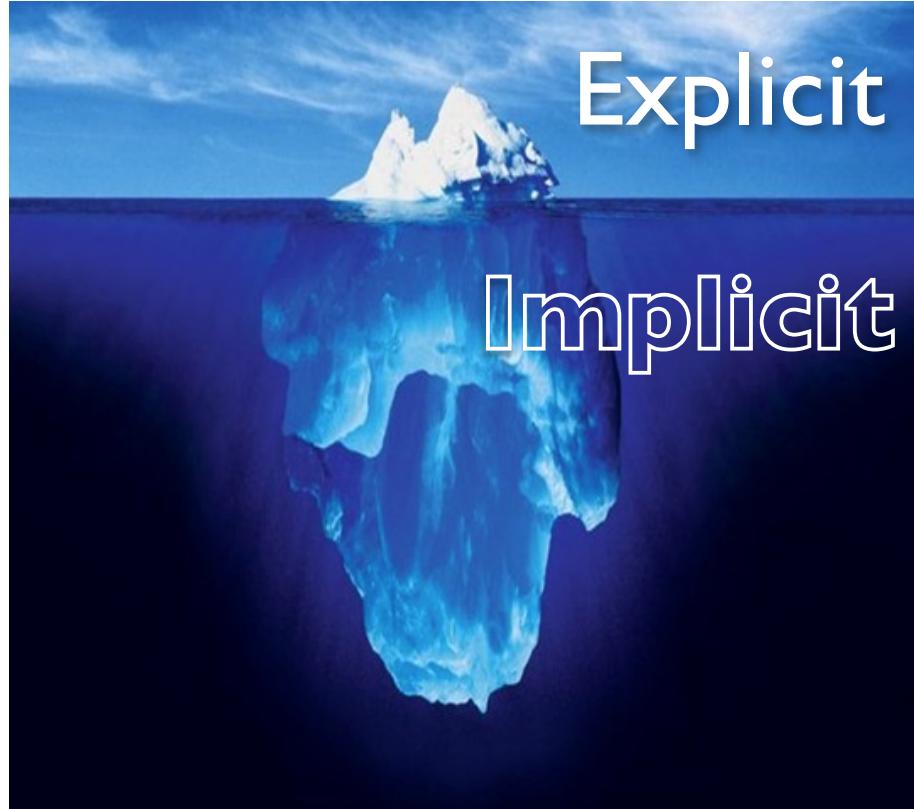
Recommenders need feedback to be useful.

Explicit

- Explicitly given
- Pro-actively acquired
- Expensive to collect

Implicit

- Indirectly given
- Larger quantity
- Latent qualities



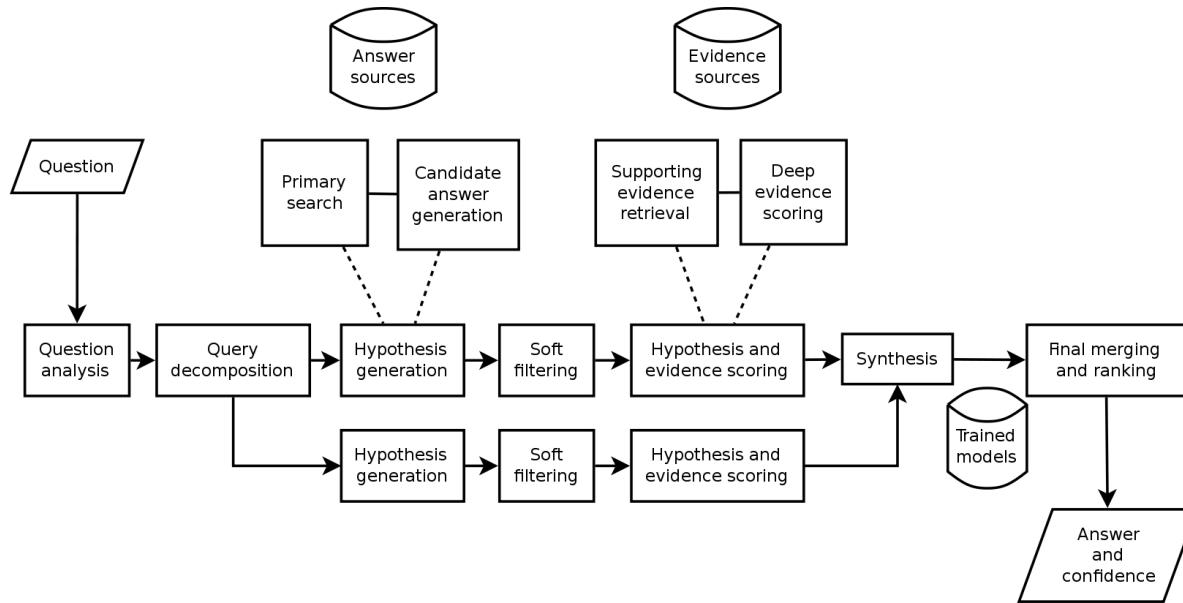
	18,000 movies					
480,000 users	x	1	1	x	...	x
	x	x	x	5	...	x
	x	x	3	x	...	x
	x	4	3	x	...	2
	...	x	x	x	...	x
	x	5	x	1	...	x
	x	x	3	3	...	x
	x	1	x	x	...	2

NOTE

This matrix will always be *sparse*!

Question Answering

IBM Watson, Wolfram Alpha



DATABASE TECHNOLOGIES: THE NOSQL MOVEMENT

Vs.

Acid

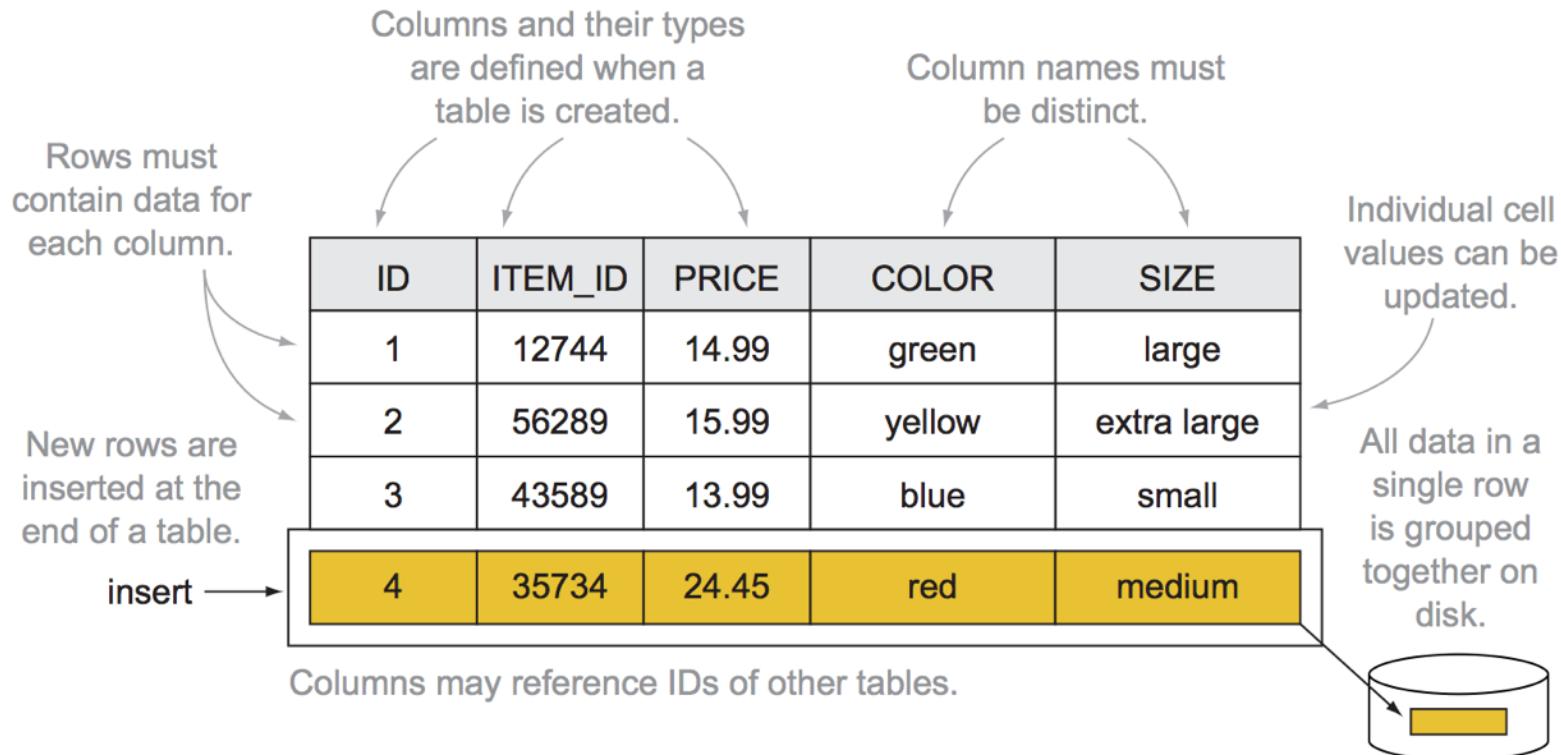
- Get transaction details right
- Block any reports while you are working
- Be pessimistic: anything might go wrong!
- Detailed testing and failure mode analysis
- Lots of locks and unlocks



Base

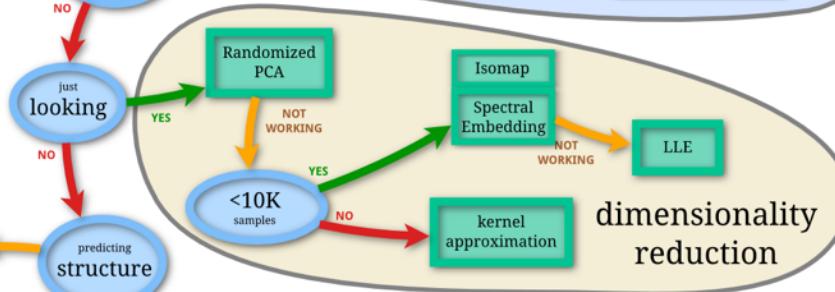
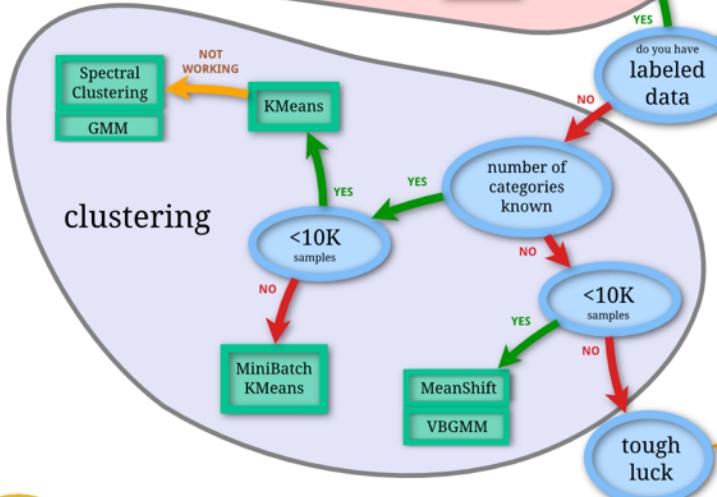
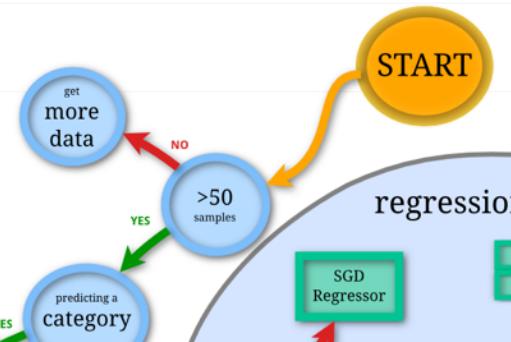
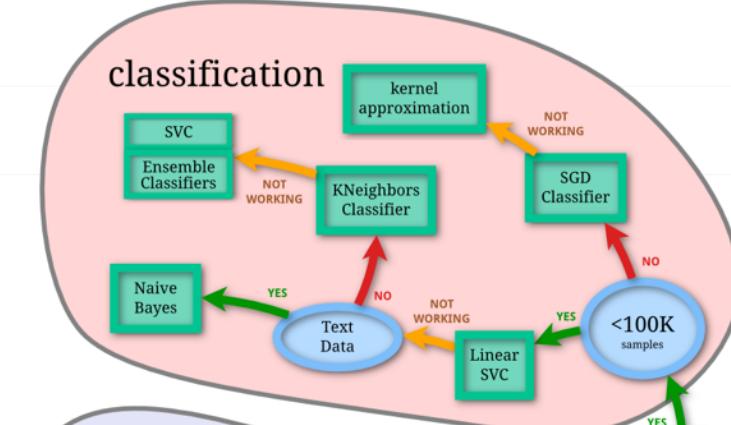
- Never block a write
- Focus on throughput, not consistency
- Be optimistic: if one service fails it will eventually get caught up
- Some reports may be inconsistent for a while, but don't worry
- Keep things simple and avoid locks





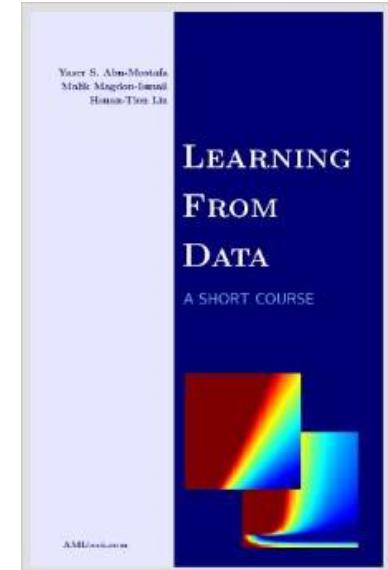
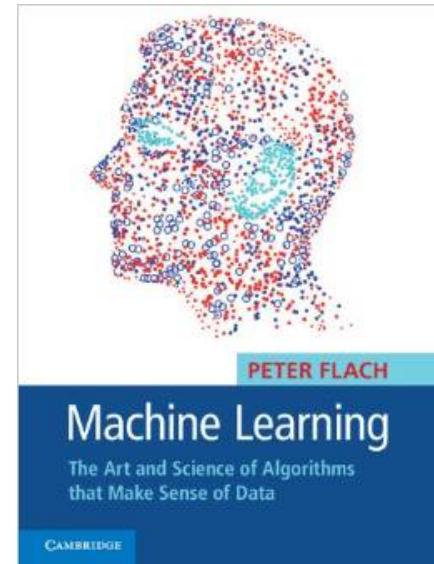
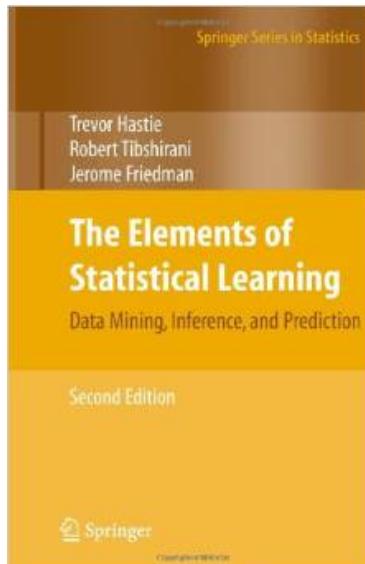
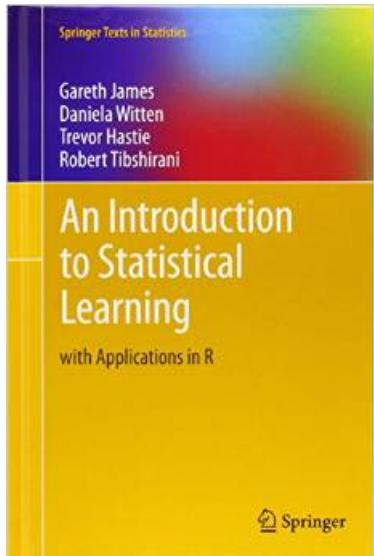
	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

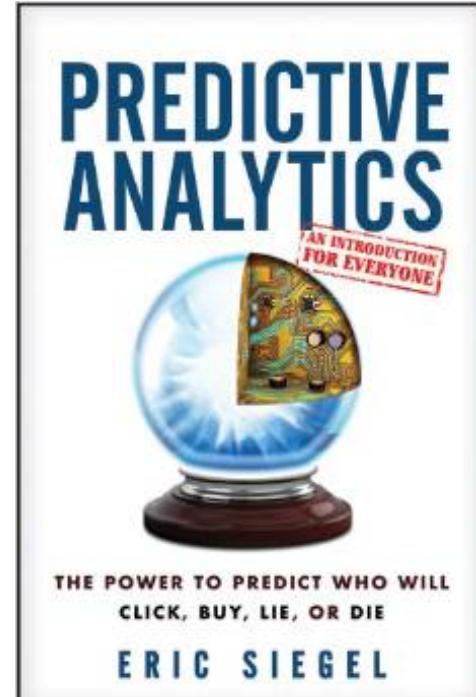
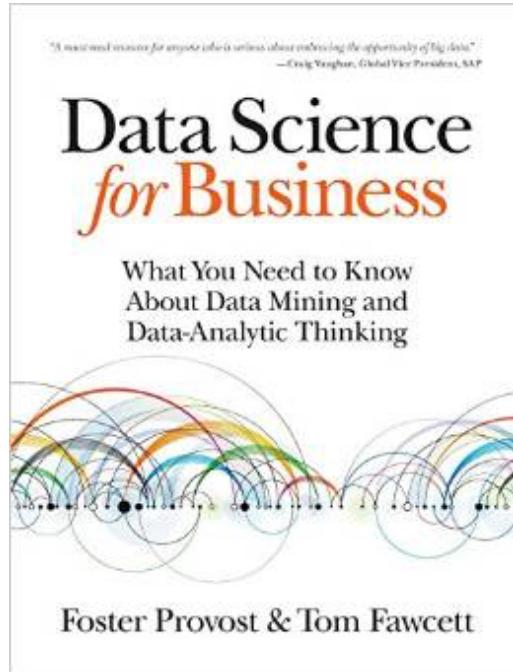
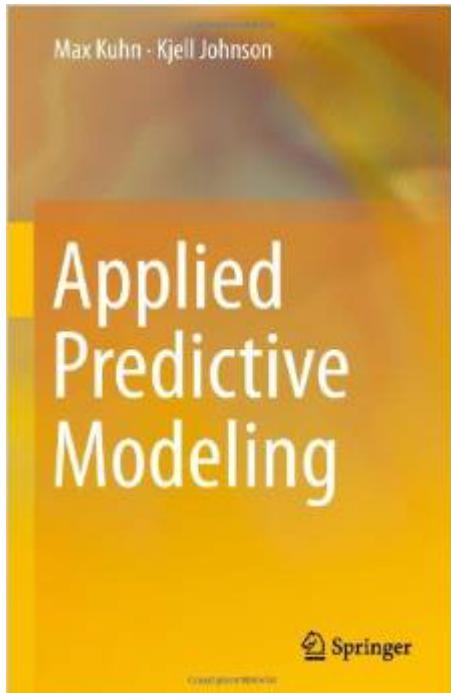
scikit-learn algorithm cheat-sheet



II. WHERE TO GO NEXT?

WHAT TO DO NEXT: READING





WHAT TO DO NEXT: READING



WHAT TO DO NEXT: ONLINE COURSES



Stanford University
Machine Learning

Ended 3 years ago

[Course Record](#)

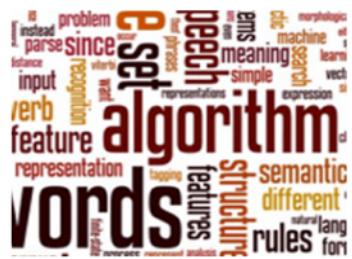


Stanford University
Mining Massive Datasets

Ended 2 months ago

[Course Record](#)

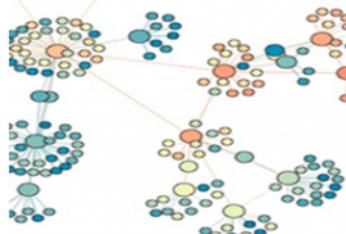
WHAT TO DO NEXT: ONLINE COURSES



Stanford University Natural Language Processing

Ended 3 years ago

[Course Record](#)



University of Michigan Social Network Analysis

Ended a year ago

[Course Record](#)

WHAT TO DO NEXT: STACKOVERFLOW PRO TIP

<http://stackoverflow.com/users/163740/ogrissel>

User Profile (Left Screenshot):

- Profile:** ogrissel (top 2% overall)
- Reputation:** 15,905
- Top Tags:** scikit-learn, python, machine-learning, svm, classification, scikits

Network Profile (Right Screenshot):

- Top Network Posts:**
 - What are common statistical sins? (75 votes)
 - Machine Learning using Python (58 votes)
 - Choice of K in K-Fold cross validation (54 votes)
 - Large scale text classification (34 votes)
 - Is there a concept of "enough" data for training statistical models? (29 votes)
 - LARS vs coordinate descent for the lasso (23 votes)
 - What is the objective Scikit-learn's Random Forest classifier is optimizing at each node? (22 votes)
- Top Posts (338):**
 - Is there a recommended package for machine learning in Python? (75 votes)
 - Save NaiveBayes classifier to disk in Scikit learn (58 votes)
 - Python: tf-idf-cosine: to find document similarity (54 votes)
 - Python scikit-learn: exporting trained classifier (34 votes)
 - Classifying Documents into Categories (29 votes)
 - Possibility to apply online algorithms on big data files with sklearn? (23 votes)
 - Is scikit-learn suitable for big data tasks? (22 votes)
 - What's the difference between LibSVM and LibLinear (21 votes)
 - How to elementwise-multiply a scipy.sparse matrix by a broadcasted dense 1d array? (18 votes)
 - fastest SVM implementation usable in python (17 votes)
- Badges (89):** GOLD, SILVER, BRONZE

- *Python: use it every day / week!*
- *Scala (if you get bored with Python; Spark has PySpark)*
- *Visualization (D3 and js)*
- *Algorithms: Artificial Neural Networks / Deep Learning*
- *“Big Data:”*
 - *MapReduce / Hadoop*
 - *Spark*

kaggle



[topcoder]TM



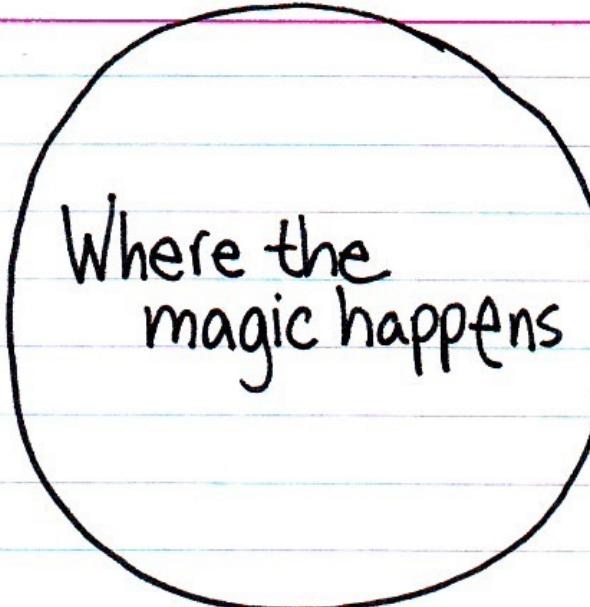
-
- *Online :*
 - <http://www.reddit.com/r/machinelearning>
 - <http://www.r-bloggers.com>
 - <http://dataelixir.com>
 - <http://www.datatau.com>
 - <http://pydata.org>
 - <http://www.technologyreview.com>
 - *LinkedIn groups on ML and DS*

Each Other!!!

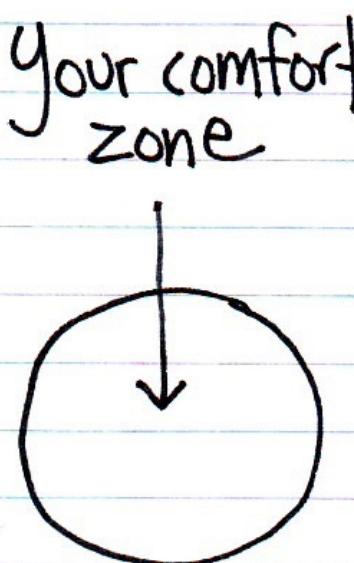
- *Class Google group / mailing list*
- *LinkedIn*
- *Study groups (Prior DAT cohorts still have dinner and share job opportunities)*

CONGRATULATIONS...

...for getting out of your comfort zone!



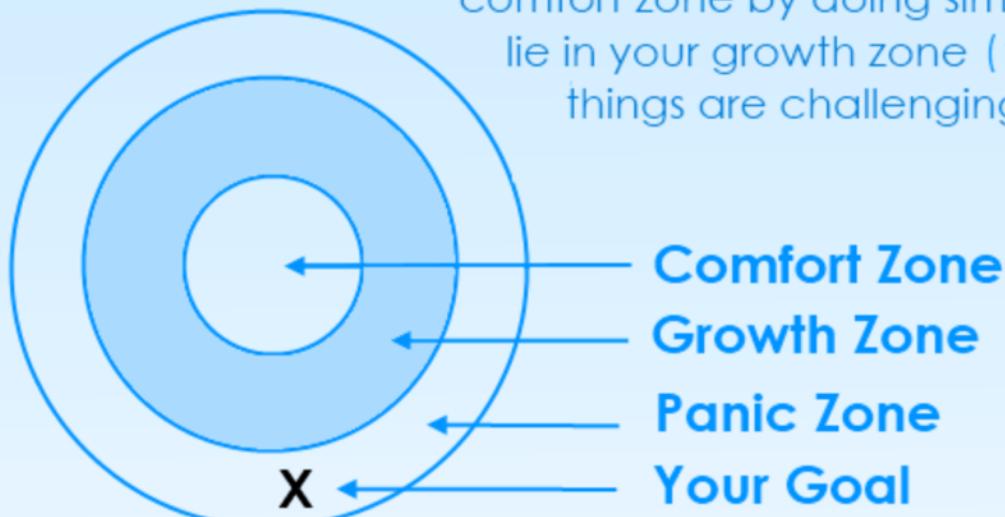
Where the
magic happens



Your comfort
zone

How to Grow Your Comfort Zone

Any goal or challenge may fall into one of three zones - your comfort zone, growth zone, or panic zone. If your goal is currently in your panic zone, i.e. it would be too scary to do now, you will need to grow your comfort zone by doing similar challenges that lie in your growth zone (the zone in which things are challenging or scary, but do-able).



How to Grow Your Comfort Zone

As you pursue challenges in your growth zone, those challenges become easier and your comfort zone expands.

Eventually, challenges that were previously in your panic zone begin to fall into your growth zone, and ultimately within your comfort zone.

