

MAPR[®] Academy

Machine Learning

Introduction to Data Science





- ▶ Features
- ▶ Algorithms and evaluation
- ▶ Over-fitting





Which is more important – the data or the algorithm?

Model = sample data + algorithm(s)



&&

Safe to assume that
everyone is using the
same algorithms



?





What is a feature?

- **Feature**
 - individual measurable property of an observation
 - Used for predicting the target (supervised learning)
- Also known as
 - *Predictor*
 - *Attribute*
 - *Variable*
- **Examples** (for a person)
 - Height (in inches)
 - Weight (in pounds)
 - Credit score (unit-less)
 - Sex (male/female)
 - Married (yes/no)
 - Nationality (American, French, Chinese, ...)
 - ZIP code (94086, ...)





Feature Categories

- **Continuous** (numeric)
 - age, time to failure, # login attempts, etc.
 - Standardize or normalize?
- **Binary**
 - yes/no, on/off, male/female.
 - Code to 0,1
- **Categorical**
 - red, yellow, green, blue
 - Also may apply to ranges of numeric values
- **Text**
 - Vocabulary (corpus dictionary)
 - N-gram (e.g. n=2 → “bigram”)

YOU are the expert in determining good features for YOUR model

Another category – words in a document, TD-IDF (would fall under transformation?)





Vectorizing Data

- Given an observation, assign features to dimensions

**Example: vectorize a text document
(bag of words approach)**

Dictionary: [a, advance, after, ..., you, yourself, youth, zigzag]



[223,1,1,0,...,12,10,6,1]

Example: vectorize an apple

Features: [size, color, weight]



[3.2, 167771848, 100]

Introduction to Machine Learning

© 2014 MapR Technologies 6

In order to cluster data, the data set must be vectorized. For example, if you wish to cluster text documents, you must define a dictionary of words and associate each word to a dimension. Such a vector might well be constituted of tens of thousands of dimensions. Given a text document (like the text representation of “Alice in Wonderland”), you’d perform a wordcount against the document to determine the vector that represents the document. Similarly, if you wish to cluster apples, you would identify a set of features that can be used to uniquely represent a given apple (for example, the radial size of the apple, the color in hex, and the weight in grams). There may be words you don’t care about that just introduce noise into the model – these are called “stop words”. Examples of stop words are those you’ll find in any document (e.g. a, an, the, some, most, all, ...).





What makes a feature “good”?

- The features are ***correlated to the target***
- Features are ***independent of each other***
- The features ***do not leak the target***
- The ***number*** of features ***is minimized***
- ***The number of features is much less than the number of data points***

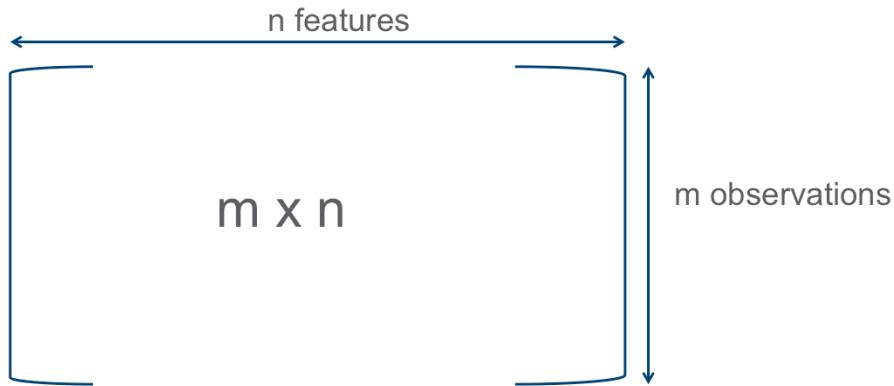
Target leak example = zip code in sioux falls (changed to fraud center)

Find targets by isolating best-performing variables, then remove & determine perf impact.





Data set dimensionality



Ideally, |features| <<< |observations|





Dimensionality reduction: Principal Component Analysis (Conceptual)

Calculate covariance matrix of original input matrix

Perform eigenvalue decomposition of covariance matrix

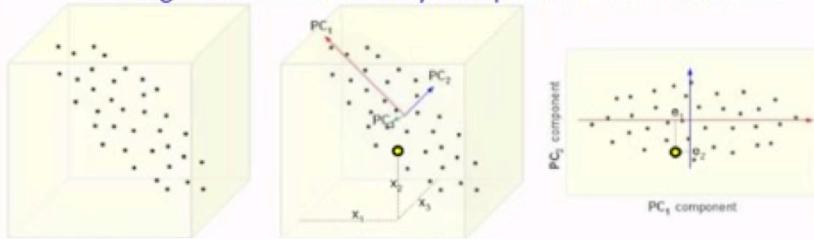
Look for orthogonal dimensions with greatest variance (principal)

Project input matrix into those dimensions (reduction)

Relative covariance is preserved across projection

Introduction to Machine

iR Technologies MAPR 9





Example: Using PCA to Reduce Dimensionality on Iris Dataset

$$iris = \begin{bmatrix} 5.1 & 3.5 & 1.4 & 0.2 \\ 4.9 & 3.0 & 1.4 & 0.2 \\ 4.7 & 3.2 & 1.3 & 0.2 \\ 4.6 & 3.1 & 1.5 & 0.2 \\ \dots & \dots & \dots & \dots \end{bmatrix} \quad \text{X} \quad pc_{1,2} = \begin{bmatrix} 0.36 & -0.66 \\ -0.85 & -0.73 \\ 0.86 & 0.17 \\ 0.36 & 0.08 \end{bmatrix} \quad \equiv \quad iris' = \begin{bmatrix} -2.68 & -0.32 \\ -2.71 & 0.18 \\ -2.89 & 0.14 \\ -2.74 & 0.32 \\ \dots & \dots \end{bmatrix}$$

$$\lambda = \begin{bmatrix} 4.2 & 0.24 & 0.08 & 0.02 \end{bmatrix} \quad \text{98%} \quad pc = \begin{bmatrix} 0.36 & -0.66 & 0.58 & 0.32 \\ -0.85 & -0.73 & -0.60 & -0.32 \\ 0.86 & 0.17 & -0.08 & -0.48 \\ 0.36 & 0.08 & -0.55 & 0.75 \end{bmatrix}$$

Introduction to Machine Learning

© 2014 MapR Technologies MAPR 10

Given iris dataset ($n=150$) with 4 features – sepal/petal height/width. Execute PCA.

Here are 4 eigenvalues and 4 principal components.

The first 2 account for 98 of total variance. So we select only the first 2 PC.

Multiplying the original matrix by PC2 matrix gives us our reconstructed data – we could use this as inputs to a model.

One drawback – if your method allows you to explain results, you will lose that since the new features are difficult to understand.





Feature scale

- **relative magnitude**

- Some features have much different magnitude than others
- Model will be dominated by largest magnitude features

- **absolute magnitude**

- Some features are extremely small and others are extremely large
- Model may lose information to underflow and overflow





Dealing with feature scale

- **Standardization**
 - transform data to measure how many standard deviations from mean
 - Example: $x' = (x - \bar{x}_{\text{mean}}) / (\sigma_{\text{SD}})$
- **Normalization:**
 - transform data such that range is between 0 and 1
 - Example: $x' = (x - \bar{x}_{\text{min}}) / (\bar{x}_{\text{max}} - \bar{x}_{\text{min}})$
- ***Both methods preserve covariance***





Feature anomalies

- **Target leaks**

- In-sample data has “extra” information not found out of sample
- (hopefully) detected during testing → model too good
- E.g. using feature “selected free shipping” to predict whether a customer will purchase or not

- **Bias**

- Sampling bias (e.g. self-selected clinical trial)
- Observer bias (e.g. measuring stick is off)
- Funding bias (e.g. fudge the data)

- **Simpson's Paradox**

- Trend that appears in different groups disappears (or reverses) when groups are combined
- E.g. average of averages != average of all points

	1995		1996		1997		Combined	
Derek Jeter	12/48	.250	183/582	.314	190/654	.291	385/1284	.300
David Justice	104/411	.253	45/140	.321	163/495	.329	312/1046	.298

Introducing machine learning

© 2014 MapR Technologies 13





- ▶ Features
- ▶ Algorithms and evaluation
- ▶ Over-fitting

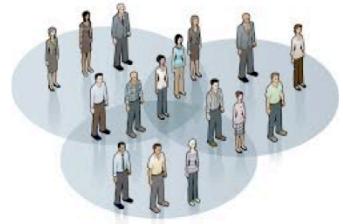


Examples of Clustering

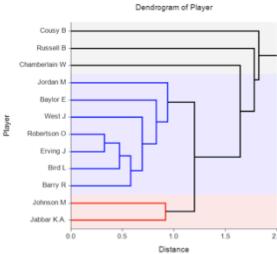
Image segmentation



Market segmentation



Basketball player style dendrogram



Introduction to Machine Learning

Phylogenetic dendrogram

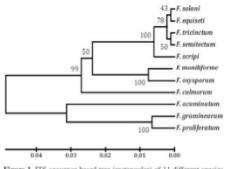


Figure 1: ITS sequence based tree (rectangular) of 11 different species of *Fusarium* constructed using UPGMA. The numbers at branch node indicate the confidence value of bootstrap replications.

© 2014 MapR Technologies  15

Clustering analysis is used in a variety of domains.

For example, in image segmentation, the objective is to identify borders of objects in an image based on the pixel data. In this example, the user has asked the app to cluster the “head” as a blue cluster and the “shirt” as a green cluster. All other data is represented in red. This could be used for computer vision.

In the case of market segmentation, the objective is to identify similarities between users based on their market habits. This could be used to target ad campaigns for specific groups of users.

You can construct dendograms using cluster analysis. A dendrogram is a hierarchical structure in which the most similar elements are closest to each other in the tree, and the most



K-Means Clustering Algorithm

- **Initialization step** assign initial coordinates to k centroids $\{m_1, m_2, \dots, m_k\}$

- **Assignment step** $S_i^{(t)} = \{x_p : \|x_p - m_i^{(t)}\|^2 \leq \|x_p - m_j^{(t)}\|^2 \forall j, 1 \leq j \leq k\},$

- **Update step** $m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$

- **Stop conditions**
 - Min change of median from last iteration
 - Sufficiently small intra-cluster distance
 - Sufficiently large intercluster distance

Introduction to Machine Learning

© 2014 MapR Technologies 16

The steps for the k-means algorithm are given above. The algorithm begins by initializing all the coordinates to centroids. There are various ways you can initialize the points, including randomly.

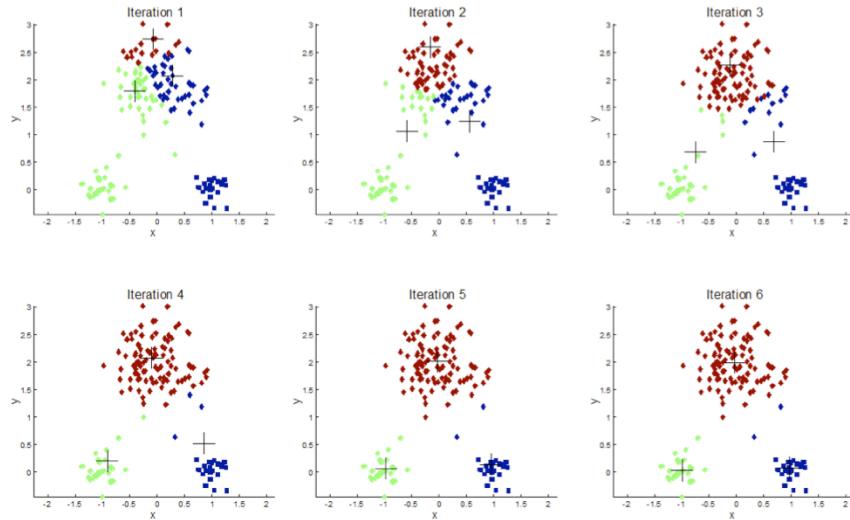
At each pass in the algorithm, you assign each point to its nearest centroid based on some distance metric (usually euclidean distance). The assignment statement reads as follows: “The ith set S at time t is equal to the set of all points x_{-p} such that the distance between x_{-p} and the mean of that set S is less than or equal to the distance between x_{-p} and the mean of all other sets S sub j where j ranges from 1 to k .”

You then update the centroids to be the “centers” of all the points assigned to it in that pass. The update step reads as





Iterations of K-Means Algorithm



Introduction to Machine Learning

© 2014 MapR Technologies MAPR 17

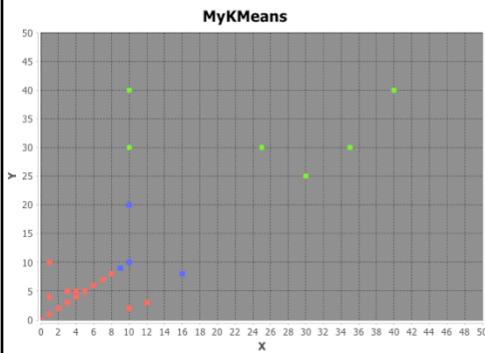
The graphic above shows 6 iterations of the k-means algorithm for a given data set. In this example, $k=3$. The centroids “move” in each iteration, and the input points stay stationary.



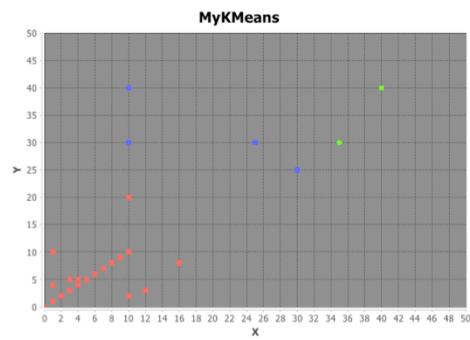


Demo: MyKMeans.java

K=3, init=random



K=3, init=elagha



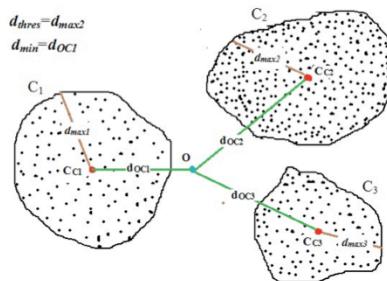
This demonstration runs the k-means algorithm using 2 different initialization methods. The first is random and the second uses an initialization method developed by El Agha et al. Note that the algorithm converges on different clusterings depending on the initialization method.





Distance-based metrics for cluster evaluation

- **Average inter-cluster distance (do I want them close or far?)**
 - Distance between centroids
 - Measures data separation between clusters
- **Average intra-cluster distance (do I want them close or far?)**
 - Distance between points in same cluster
 - Measures data similarity within clusters



In order to evaluate an algorithm, you have to have metrics with which to evaluate. These are either minimization or maximization metrics. For example, you can consider the average distance between centroids as the metric you wish to optimize. Depending on your use case, you may wish to minimize or maximize this distance. You may also consider the average distance between points and their centroid as the metric to optimize. Again, depending on your use case, you may wish to minimize or maximize this distance.





What is Naïve Bayes?

NB assumes:

- class C is dependent on features F_1, F_2, \dots, F_n
- all the features F_1, F_2, \dots, F_n are independent

$$p(C|F_1, \dots, F_n) = \frac{p(C) p(F_1, \dots, F_n|C)}{p(F_1, \dots, F_n)}.$$

the “naïvete” comes from the independence assumption

Naïve Bayes is a machine learning algorithm that is based on Bayes theorem. The crux of Naïve Bayes (the “naïve” part) is based on the assumption that all the features are independent of each other. Even though this may be too strong an assumption, Naïve Bayes has been shown to work well in several domains.





Demo: Use NB to build a simple, basic, naïve spam filter

1. Calculate probability that the email is spam

quotient of (number of positives) to (number of positives + number of negatives)

product of positive frequencies of terms in subject line

$$P(\text{spam}|\text{subject-line}) = P(\text{spam}) * P(\text{subject-line}|\text{spam}) / P(\text{subject-line})$$

2. Calculate probability that the email is ham

quotient of (number of negatives) to (number of positives + number of negatives)

product of negative frequencies of terms in subject line

$$P(\text{ham}|\text{subject-line}) = P(\text{ham}) * P(\text{subject-line}|\text{ham}) / P(\text{subject-line})$$

3. Determine which is more likely

if $P(\text{spam}|\text{subject-line}) > P(\text{ham}|\text{subject-line})$ return "spam"
else return "ham"

In this demo, we train and query a Naïve Bayes classifier for detecting spam email. This particular algorithm only considers the contents of the subject line in the model.





Demo: MySpamClassifier.java

```
$ java MySpamClassifier DATA/train.txt
num emails is 8717
enter subject line: send your password
probability of not spam is 8.849381E-24
normalized probability of not spam is 0.12610133
probability of spam is 6.1327366E-23
normalized probability of spam is 0.8738987
==> model predicts spam
=====
enter subject line: vote now!
probability of not spam is 2.0945408E-6
normalized probability of not spam is 0.4538798
probability of spam is 2.520207E-6
normalized probability of spam is 0.5461202
==> model predicts spam
=====
enter subject line: vote now
probability of not spam is 8.105502E-5
normalized probability of not spam is 0.89094615
probability of spam is 9.921322E-6
normalized probability of spam is 0.10905387
==> model predicts not spam

cat test-ham-100.txt | java MySpamClassifier DATA/train.txt | grep -c "model predicts not spam"
84

cat test-spam-100.txt | java MySpamClassifier DATA/train.txt | grep -c "model predicts spam"
75
```

Introduction to Machine Learning

© 2014 MapR Technologies **MAPR** 22

This particular NB program takes as input a line of text which represents the subject line of an email and outputs information regarding whether the model considers it spam or not. Interestingly, the string “vote now” is not considered spam while the same subject line with an exclamation point is considered spam. This is due to the fact that many of the spam emails in the training corpus contain the exclamation mark in their subject line.





Measures based on TP/FP/TN/FN

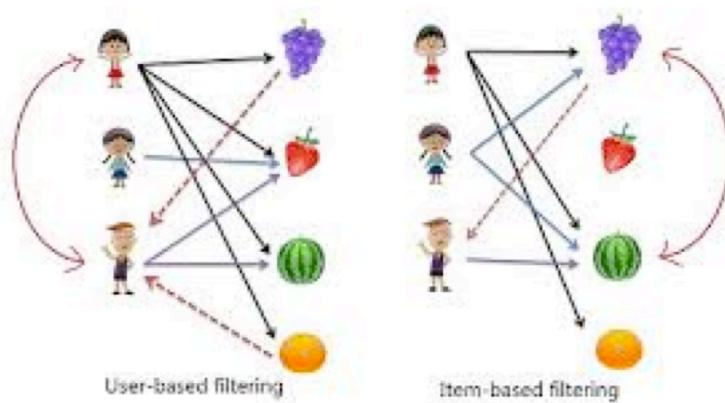
- **Sensitivity** (aka true positive rate)
$$\text{TPR} = \text{TP} / (\text{TP} + \text{FN})$$
- **Specificity** (aka true negative rate)
$$\text{TNR} = \text{TN} / (\text{TN} + \text{FP})$$
- **Accuracy**
$$\text{ACC} = (\text{TP} + \text{TN}) / (\text{P} + \text{N})$$
- Different domains have different “penalties”
e.g. spam detection vs cancer detection

Based on TP, FP, TN, and FN, you can calculate rate metrics to evaluate your classifier, as shown in the slide above. The higher these values, the better the model.





Recommendation: User-based vs item-based filtering



Introduction to Machine Learning

© 2014 MapR Technologies **MAPR** 24

Users are associated with item ***preferences*** through history

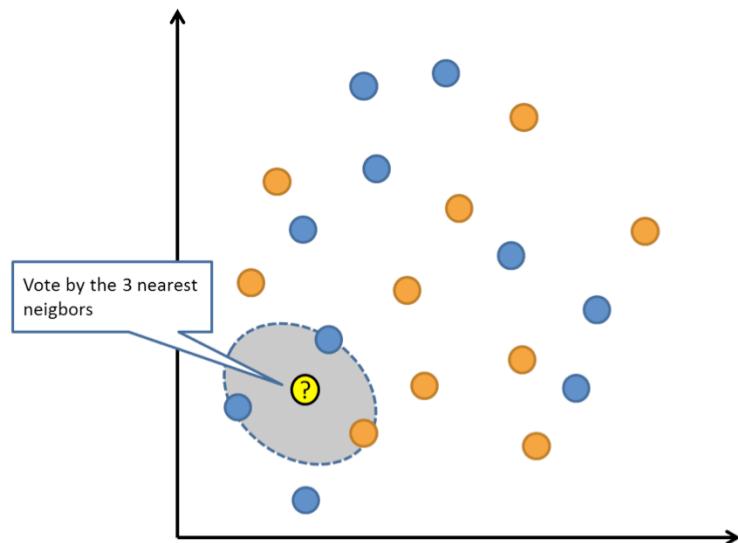
Similarity is constructed between users (or items)

Recommendations based on ***similarity to other users (or items)***





Recommendation: K-Nearest Neighbor (KNN)



Introduction to Machine Learning

© 2014 MapR Technologies **MAPR** 25



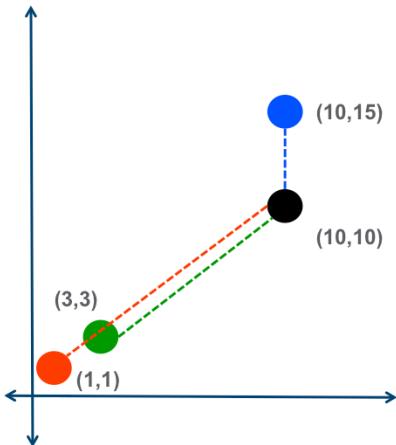
Recommendation: KNN algorithm

- **Distance**
 - Euclidean
 - Manhattan
 - Cosine
- **Majority voting**
 - Problematic if class distribution is skewed
 - Could use distance as weight to compensate
 - Pick odd number for k to avoid ties





Recommendation Demo: MyKNearestNeighbor.java



```
java MyKNearestNeighbor 1 10 10 euclidean
input: [10 10]
[1 1]
distance is 12.727922061357855
[3 3]
distance is 9.899494936611665
[10 15]
distance is 5.0
k nearest neighbors are:
[10 15]
```

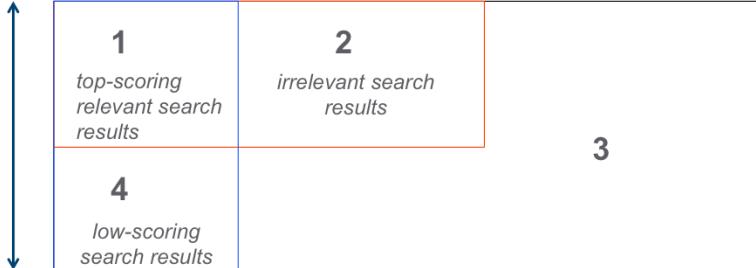
```
java MyKNearestNeighbor 1 10 10 cosine
input: [10 10]
[1 1]
distance is 0.9999999999999998
[3 3]
distance is 1.0
[10 15]
distance is 0.9805806756909201
k nearest neighbors are:
[3 3]
```

 Recommendation: Precision and recall in document search

All docs in corpus	= 1 + 2 + 3 + 4
All results from search	= 1 + 2
Top results from search	= 1
Relevant documents	= 1 + 4

$\text{precision} = |1| / (|1| + |2|)$

$\text{recall} = |1| / (|1| + |4|)$



The diagram shows a 2x2 grid representing search results:

- Top-left cell (blue border): **1** top-scoring relevant search results
- Top-right cell (red border): **2** irrelevant search results
- Bottom-left cell (blue border): **4** low-scoring search results
- Bottom-right cell (no border): **3** (empty)

Precision and recall are used to **evaluate search results**

Precision = proportion of top-scoring results that are relevant

Recall = proportion of relevant results that are top-scoring





- ▶ Features
- ▶ Algorithms and evaluation
- ▶ Over-fitting





Overfitting and underfitting





Dealing with overfitting

- **Regularization**

- Occam's razor (William of Occam): most plausible explanation is the simplest
- Impose penalty in cost function based on parameter complexity

$$\min(E_{\text{out}}(X, Y)) \rightarrow \min(E_{\text{in}}(X, Y) + \lambda \| w \|)$$

Smaller $\lambda \rightarrow$ overfitting; larger $\lambda \rightarrow$ underfitting

- **Smoothing**

- Sunrise problem (Laplace): $P(\text{sun will rise tomorrow}) = ?$
- Use probability distributions for parameter estimates for sparse data

$$\hat{\theta}_i = \frac{x_i + \alpha}{N + \alpha d} \quad (i = 1, \dots, d),$$

Smaller $\alpha \rightarrow$ overfitting; larger $\alpha \rightarrow$ underfitting





Dealing with overfitting: cross-validation

- Leave k out
- Train on N-k points, validate on k points; gives estimate of out-of-sample error (which we want to minimize)
- Example (leave 1 out)
 - Take (x_1, y_1) out, train on N-1; calculate e_1
 - Put (x_1, y_1) back in training set
 - Take (x_2, y_2) out, train on N-1; calculate e_2
 - Put (x_2, y_2) back in training set
 - ...
 - Take (x_N, y_N) out, train on N-1; calculate e_N
 - Calculate cross-validation error = $1/N * (e_1 + e_2 + \dots + e_N)$
- Repeat above steps, each for a *different model* (e.g. higher-order polynomial, different λ , different α , ...)
- select model with lowest average cross-validation error





K-fold cross-validation

- Leave 1 out → lots of iterations per epoch 😞
- Cross-validation generalizes to higher values of k (k-fold cross-validation)
- $k=N/10$ (10-fold cross-validation) is empirically a good rule of thumb
- Leave first k out, train on $N-k$, calculate e_{k1} ; repeat 10 times; calculate average e ; perturb model and repeat cross-validation k times; choose model
- Question: what is the difference between validation and testing?

