



Computer
Science

CSC 480/580: Principles of Machine Learning

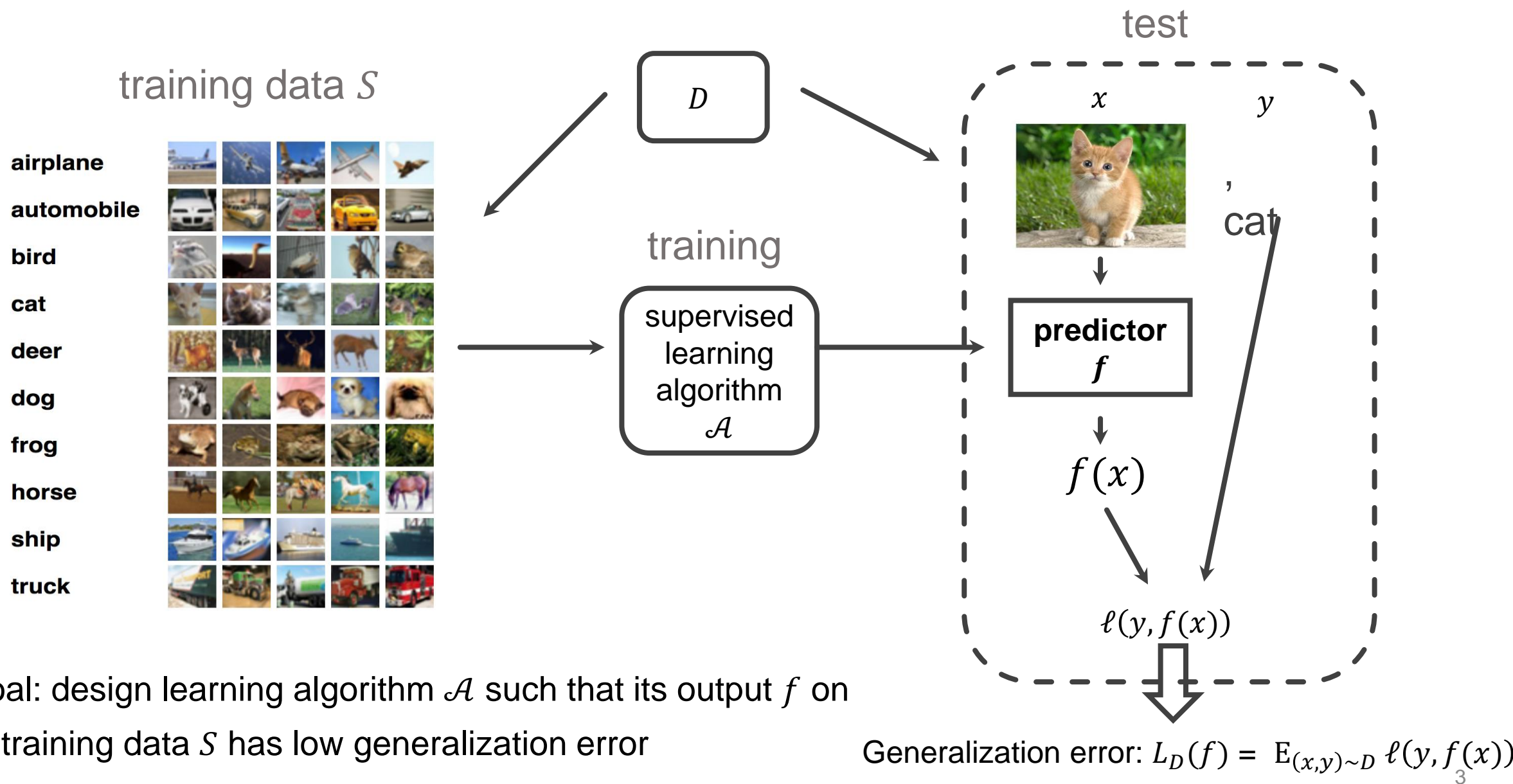
Midterm Review

Chicheng Zhang

Midterm

- Mar 12, in class (9:30-10:45am)
- About 6 questions
- You can bring a A4-size note
 - The process of preparing for the note is a good way to organize your knowledge

Supervised learning setup: putting it together

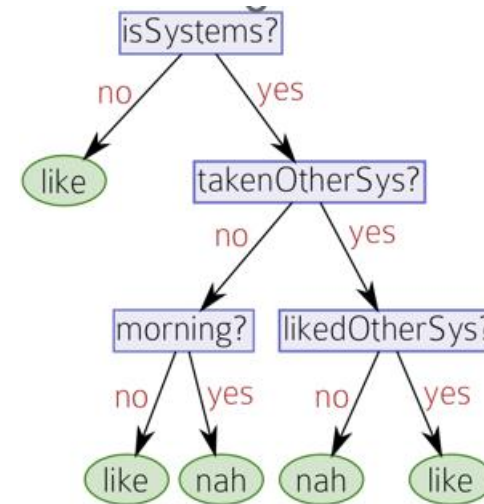


Decision trees

- Test: predict using a decision tree:

Algorithm 2 `DECISIONTREETEST`(*tree*, *test point*)

```
1: if tree is of the form LEAF(guess) then  
2:   return guess  
3: else if tree is of the form NODE(f, left, right) then  
4:   if f = no in test point then  
5:     return DECISIONTREETEST(left, test point)  
6:   else  
7:     return DECISIONTREETEST(right, test point)  
8:   end if  
9: end if
```



guess=prediction

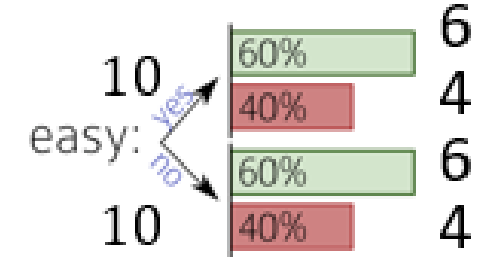
left=no
right=yes

- Training: how to design a learning algorithm \mathcal{A} that can build trees f from training data?

Algorithm 1 `DECISIONTREETRAIN`(*data*, *remaining features*)

```
1: guess ← most frequent answer in data // default answer for this data
2: if the labels in data are unambiguous then
3:   return LEAF(guess) // base case: no need to split further
4: else if remaining features is empty then
5:   return LEAF(guess) // base case: cannot split further
6: else // we need to query more features
7:   for all  $f \in \textit{remaining features}$  do
8:     NO ← the subset of data on which  $f=no$ 
9:     YES ← the subset of data on which  $f=yes$ 
10:    score[ $f$ ] ← # of majority vote answers in NO
11:                + # of majority vote answers in YES
12:   end for
13:    $f$  ← the feature with maximal score( $f$ )
14:   NO ← the subset of data on which  $f=no$ 
15:   YES ← the subset of data on which  $f=yes$ 
16:   left ← DECISIONTREETRAIN(NO, remaining features \ { $f$ })
17:   right ← DECISIONTREETRAIN(YES, remaining features \ { $f$ })
18:   return NODE( $f$ , left, right)
19: end if
```

answer=label
unambiguous=achieves 100% acc



$\text{Score}(f, S) =$ “informativeness of f (in predicting y) for dataset S ”

Q: is this algorithm guaranteed to terminate?

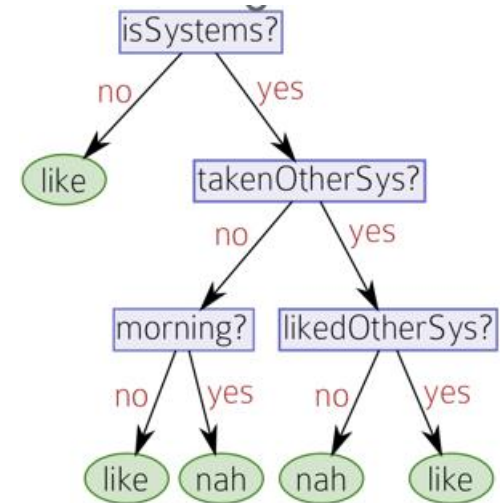
Example questions

- Given a decision tree and a test example, find out the tree's prediction on the example

- Calculate the information score of a feature

$$\text{Score}(S, f) = u(S) - (p_L u(S_L) + p_R u(S_R))$$

- See. E.g. HW1, Problem 2



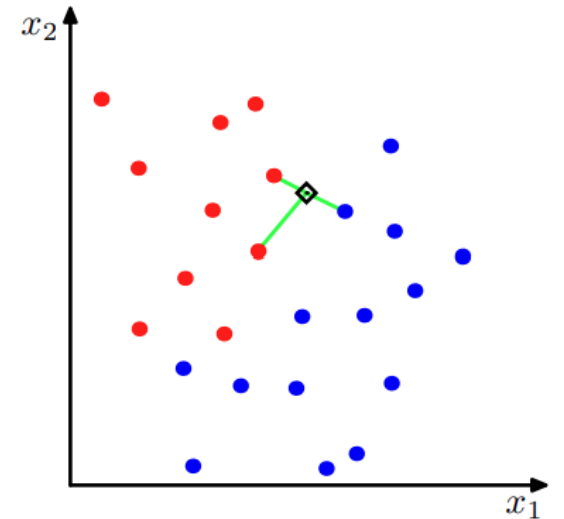
k -nearest neighbors (k -NN): main concept

Training set: $S = \{ (x_1, y_1), \dots, (x_m, y_m) \}$

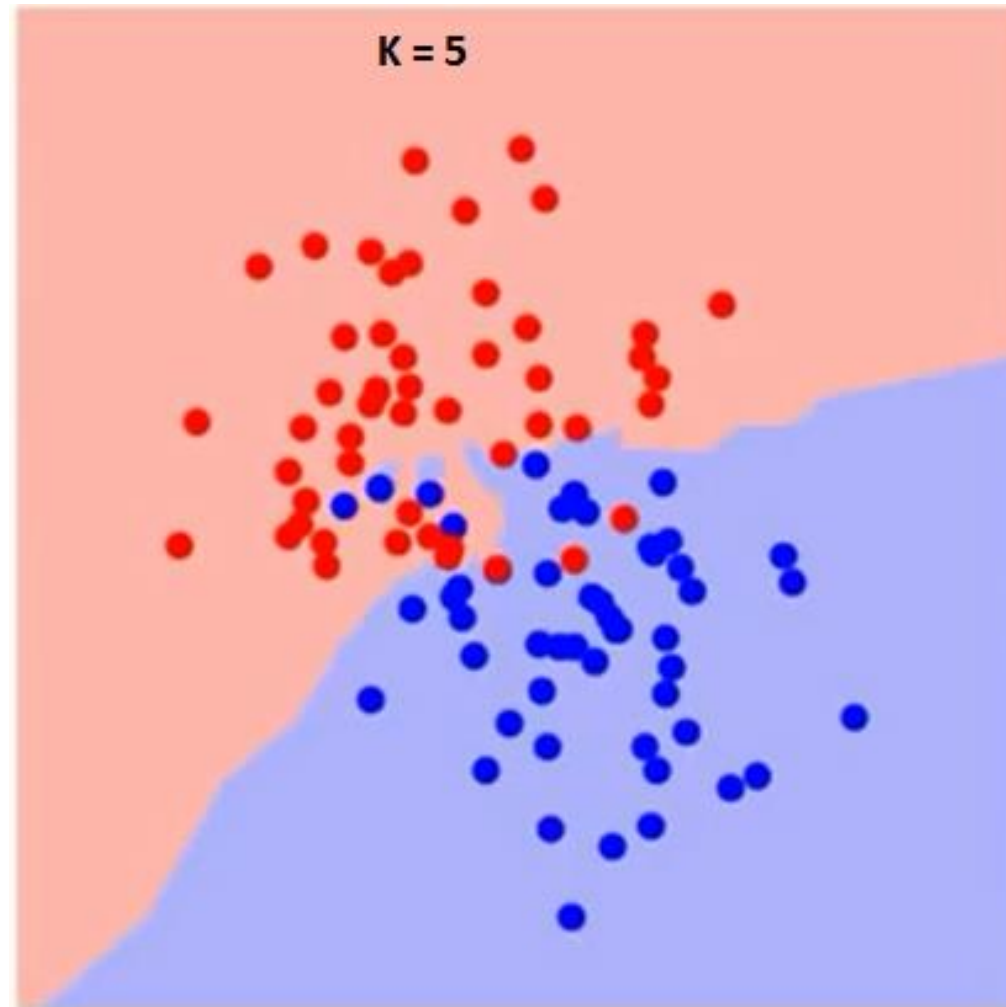
Inductive bias: given test example x , its label should resemble the labels of **nearby points**

Function

- input: x
- find the k nearest points to x from S ; call their indices $N(x)$
- output: the majority vote of $\{y_i : i \in N(x)\}$
 - For regression, the average.



k-NN classification example



← decision boundary

Variations

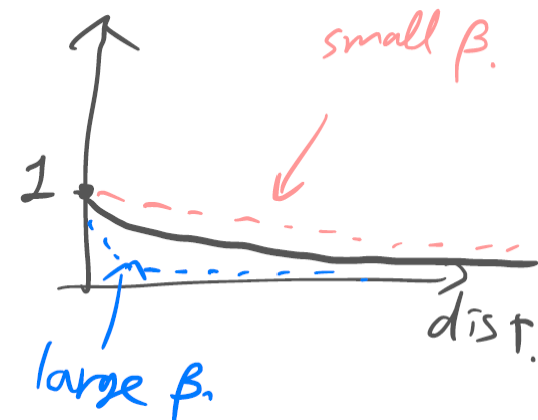
- Classification

- Recall the majority vote rule: $\hat{y} = \arg \max_{y \in \{1, \dots, C\}} \sum_{i \in N(x)} 1\{y_i = y\}$

- Soft weighting nearest neighbors: $\hat{y} = \arg \max_{y \in \{1, \dots, C\}} \sum_{i=1}^m w_i 1\{y_i = y\}$,

where $w_i \propto \exp(-\beta d(x, x_i))$, or $\propto \frac{1}{1+d(x, x_i)^\beta}$

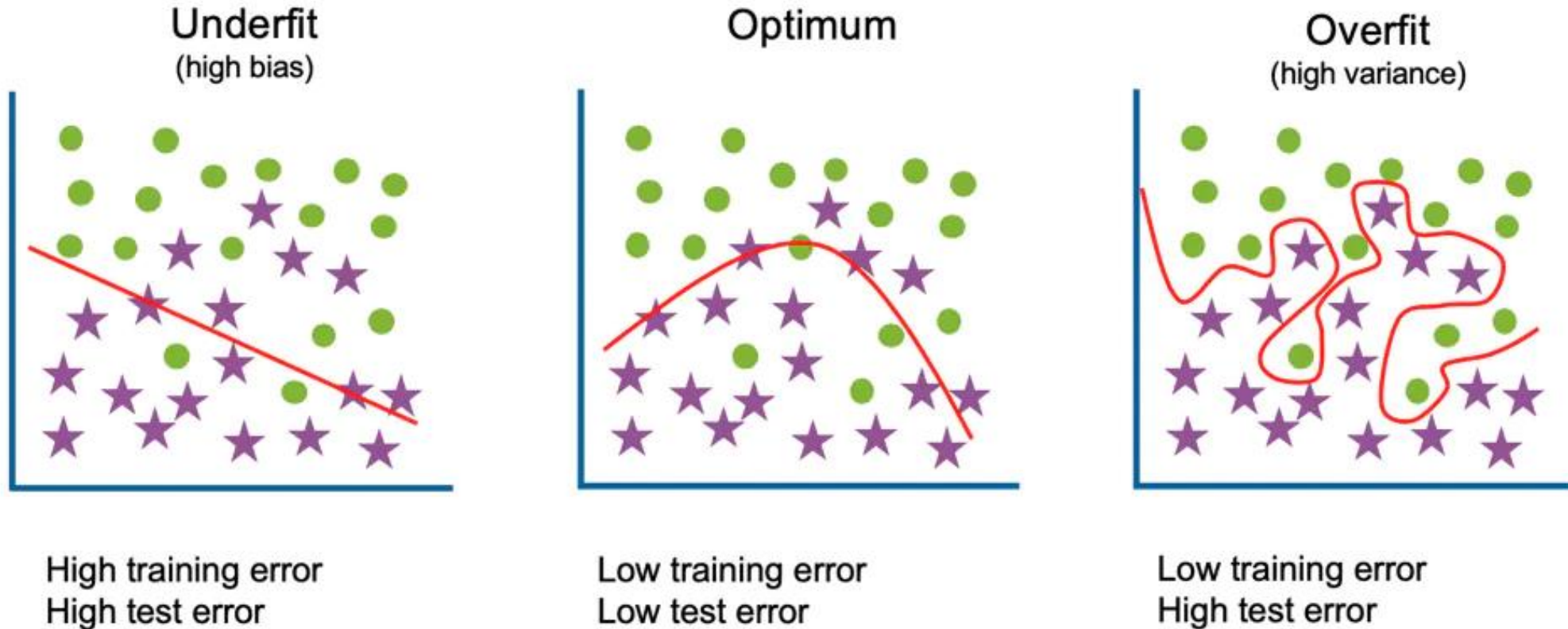
- Quick question: is a larger β resembling a larger or smaller k ?



Example question

- Given a dataset (say, of size 5), compute its induced k -nearest neighbor prediction on a test example
- Given a training dataset, compute the training error of its induced 3-nearest neighbor classifier

Overfitting vs Underfitting



Example question

- Give three scenarios of underfitting / overfitting when training a classifier
 - E.g. k-nearest neighbor, decision tree, logistic regression, training linear classifier using Perceptron
- How can we combat underfitting / overfitting?

Bayes optimal classifier

Theorem f_{BO} achieves the smallest 0-1 error among all classifiers.

$$f_{BO}(x) = \arg \max_{y \in \mathcal{Y}} P_D(X = x, Y = y) = \arg \max_{y \in \mathcal{Y}} P_D(Y = y | X = x), \forall x \in \mathcal{X}$$

Example Iris dataset classification:



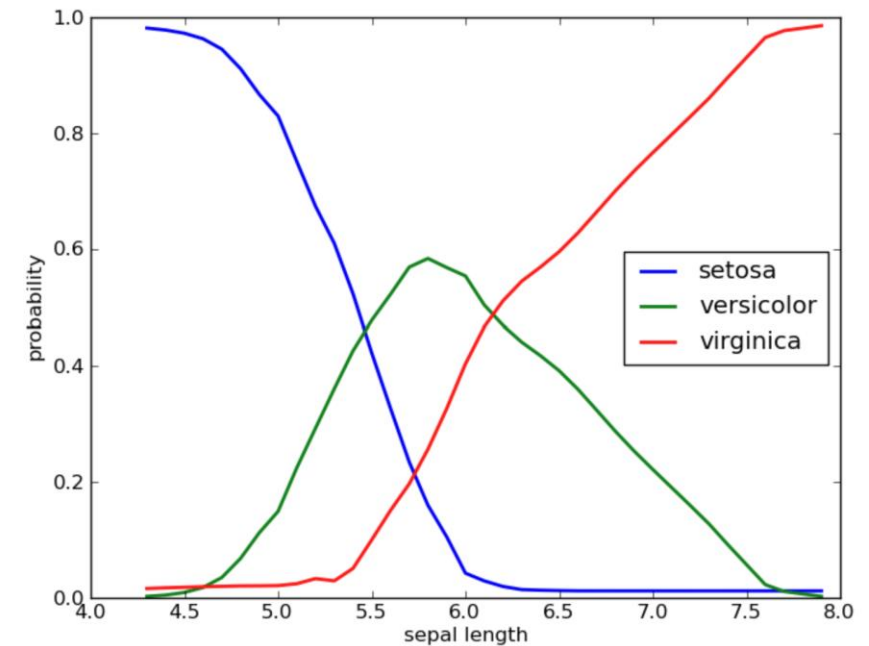
Iris Setosa



Iris Versicolor



Iris Virginica

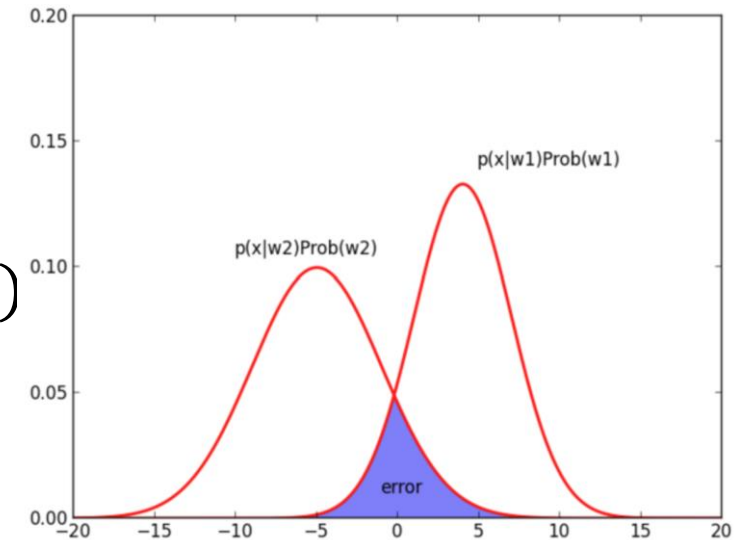
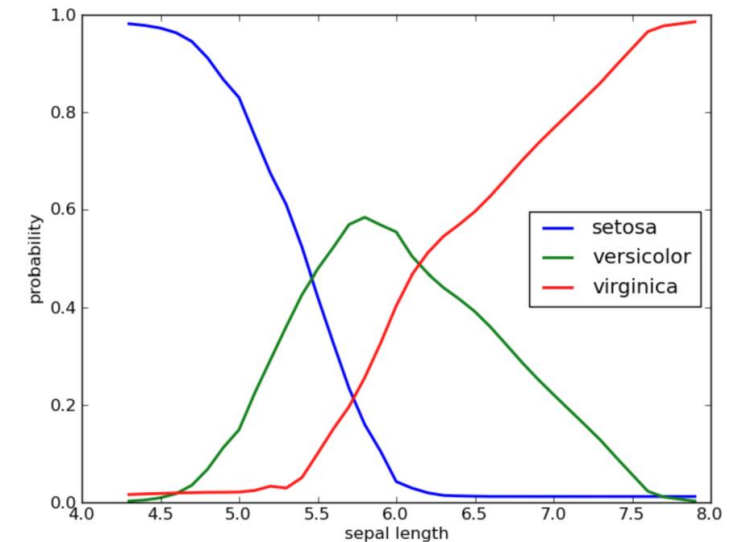


Bayes error rate: alternative form

$$\begin{aligned} L_D(f_{BO}) &= \sum_x \left(1 - \max_y P_D(Y = y | X = x) \right) P_D(X = x) \\ &= E \left[1 - \max_y P_D(Y = y | X) \right] \end{aligned}$$

- Special case: binary classification

$$L_D(f_{BO}) = \sum_x \min(P_D(Y = +1, X = x), P_D(Y = -1, X = x))$$

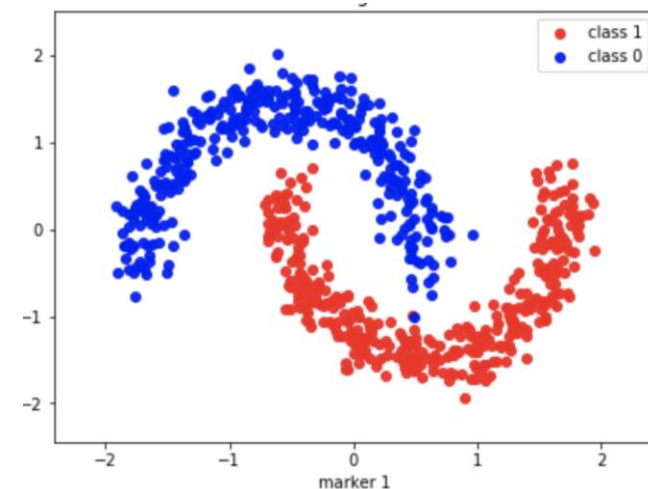
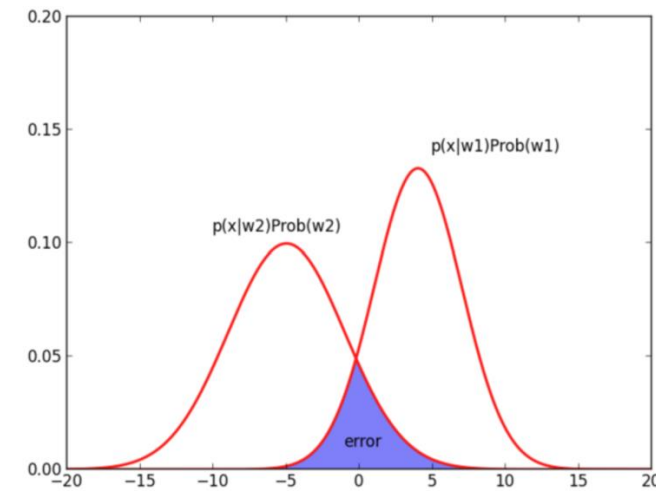
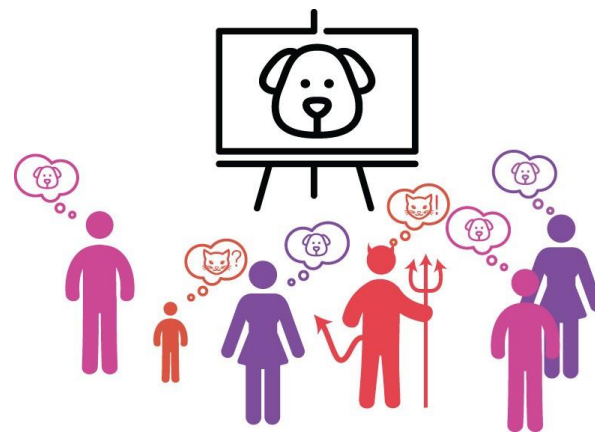


Example question

$$L_D(f_{BO}) = \sum_x \min(P_D(Y = +1, X = x), P_D(Y = -1, X = x))$$

• Give two real-world example data distributions whose Bayes errors are nonzero

- Limited feature representation
- Noise in the training data
- ...
- (see slides)



Example question

- Given a data distribution D over (x, y) , compute the error rate of a classifier & D 's Bayes error rate
 - See HW1, problem 1

New measures of classification performance

- True positive rate (TPR)

$$= \frac{TP}{P} = \frac{P(\hat{y}=+1, y=+1)}{P(y=+1)}$$

(aka recall, sensitivity)

- True negative rate (TNR) = $\frac{TN}{N}$

(specificity)

- False positive rate (FPR) = $\frac{FP}{N}$

- False negative rate (FNR) = $\frac{FN}{P}$

- Precision = $\frac{TP}{P\text{-called}} = \frac{P(\hat{y}=+1, y=+1)}{P(\hat{y}=+1)}$, P - called = TP + FP

The diagram illustrates a confusion matrix. The horizontal axis is labeled 'actual class' and is divided into 'positive' and 'negative'. The vertical axis is labeled 'predicted class' and is divided into 'positive' and 'negative'. The matrix cells are: top-left (positive predicted, positive actual) is 'true positives (TP)'; top-right (positive predicted, negative actual) is 'false positives (FP)' with a blue label 'Type I error' below it; bottom-left (negative predicted, positive actual) is 'false negatives (FN)' with a blue label 'Type II error' below it; bottom-right (negative predicted, negative actual) is 'true negatives (TN)'. Red brackets group the columns under 'actual class' and the rows under 'predicted class'.

		actual class	
		positive	negative
predicted class	positive	true positives (TP)	false positives (FP) Type I error
	negative	false negatives (FN) Type II error	true negatives (TN)

$$P = TP + FN \quad N = FP + TN$$

Example question

- Given a small binary classification dataset (say of size 5) and a predictor's prediction score on it; compute the ROC curve
- What are the x and y values in the ROC curve?
- X: FPR
- Y: TPR

$c(x_i)$	y_i
.99	+
.98	+
.72	-
.51	-
.24	+

Linear Regression

Regression Learn a function that predicts outputs from inputs,

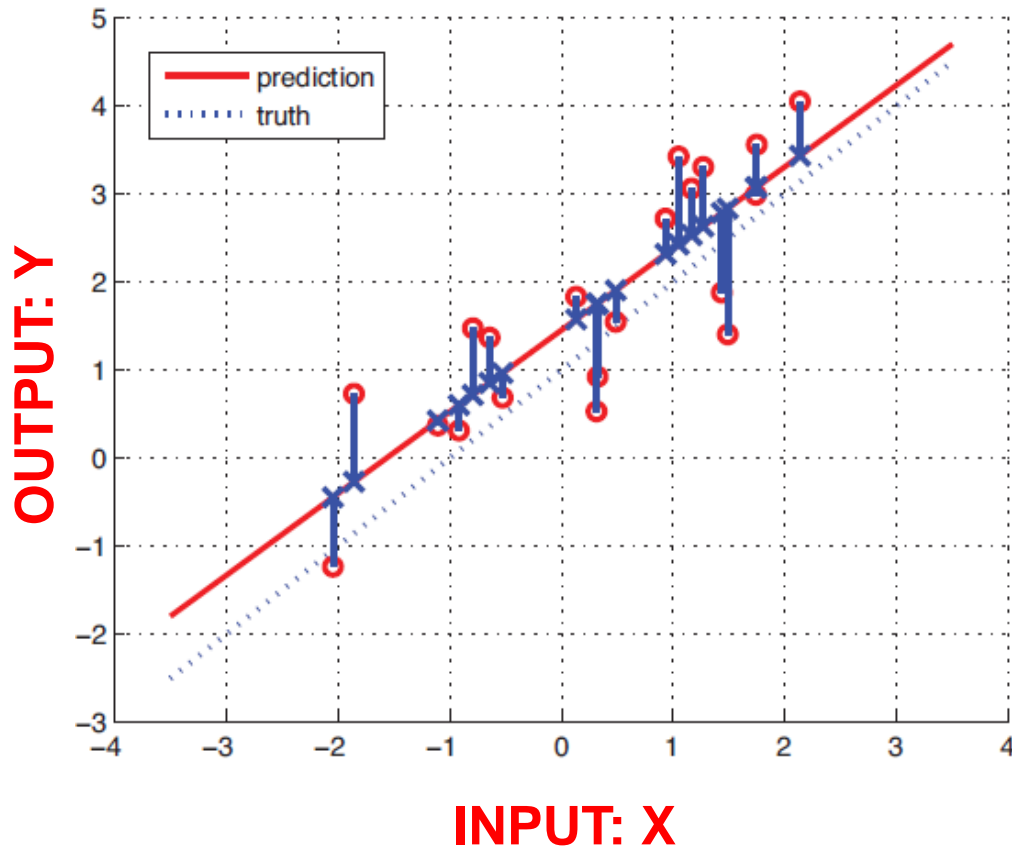
$$y = f(x)$$

Outputs y are real-valued

Linear Regression As the name suggests, uses a *linear function*:

$$y = w^T x + b$$

We will add noise later...



Linear Regression

Input-output mapping is not exact, so we will add zero-mean Gaussian noise,

$$y = w^T x + \epsilon \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

**Multivariate Normal
(uncorrelated)**

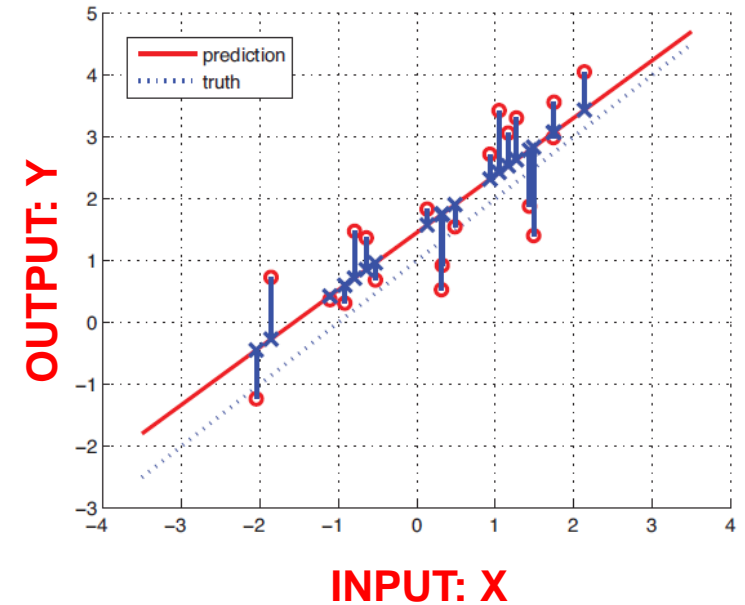
This is equivalent to the likelihood function,

$$p(y | w, x) = \mathcal{N}(y | w^T x, \sigma^2)$$

Because Adding a constant to a Normal RV is still a Normal RV,

$$z \sim \mathcal{N}(m, P) \quad z + c \sim \mathcal{N}(m + c, P)$$

In the case of linear regression $z \rightarrow \epsilon$ and $c \rightarrow w^T x$



Learning linear regression models

We need to learn the model from data
by learning the regression weights

Data – We have this

$$y = w^T x + \epsilon$$

Random; Can't do
anything about it

Don't know these;
need to learn them

How to do this?
What makes *good*
weights?

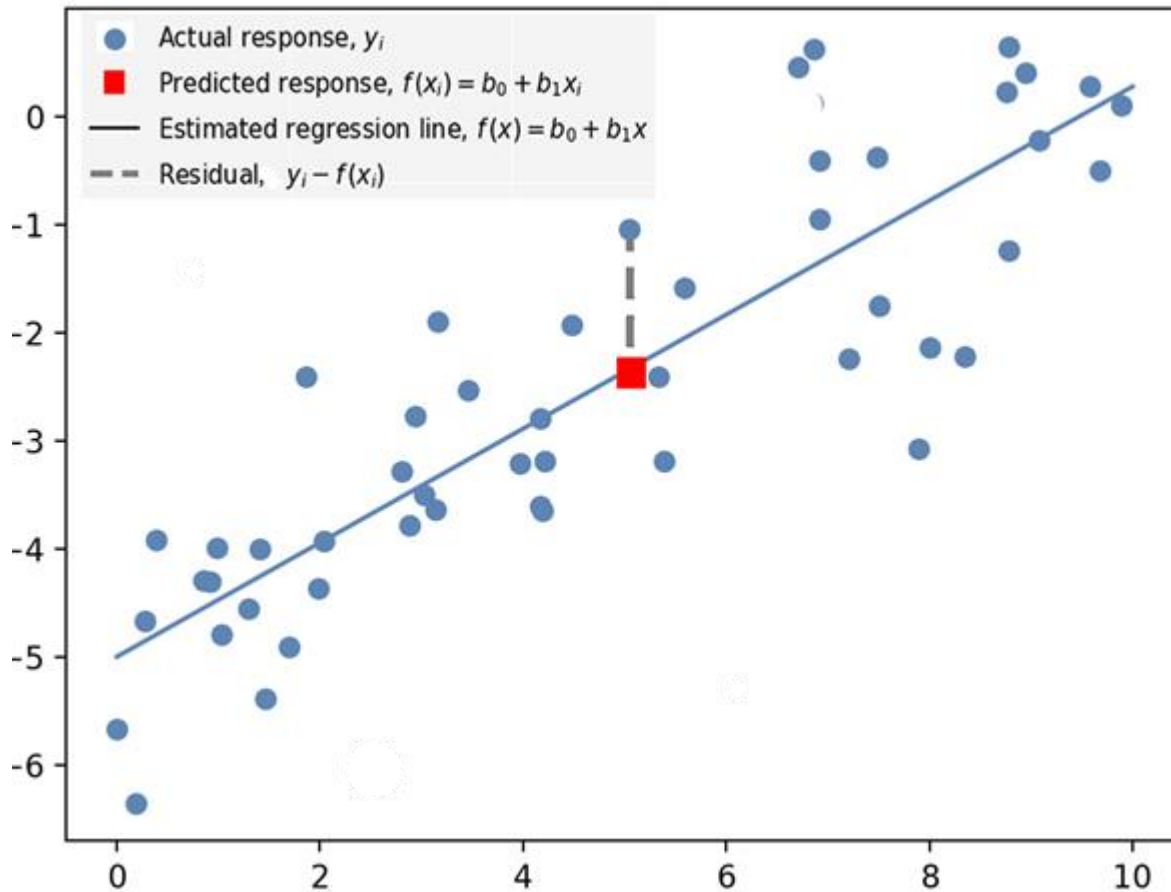
Learning Linear Regression Models

There are several ways to think about fitting regression:

- **Intuitive** Find a plane/line that is close to data
- **Functional** Find a line that minimizes the *least squares* loss
- **Estimation** Find maximum likelihood estimate of parameters

They are all the same thing...

MLE for Linear Regression



Substitute linear regression prediction into MLE solution and we have,

$$\min_w \sum_{i=1}^N (y_i - wx_i)^2$$

So for Linear Regression,
MLE = Least Squares
Estimation

MLE of Linear Regression

Using previous results, MLE is equivalent to minimizing squared residuals,

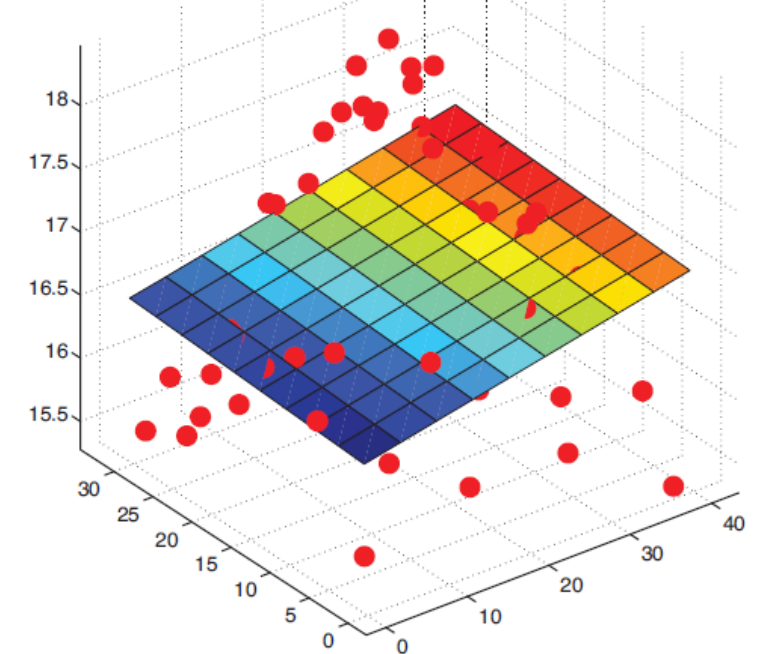
$$\min_w \sum_{i=1}^N (y_i - w^T x_i)^2 = \|\mathbf{y} - w^T \mathbf{X}\|^2$$

Some slightly more advanced linear algebra gives us a solution,

$$w = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Ordinary Least Squares (OLS) solution

[Image: Murphy, K. (2012)]



- Derivation a bit involved for lecture but...
- We know it has a closed-form and why
 - We can evaluate it
 - Generally know where it comes from

Example question

- Given a small dataset (say, size 5):
 - Write down the optimization problem least square regression solves
 - Be able to write down the derivative of the objective wrt w
 - Compute the least square linear regression coefficient

- Given a small dataset, write down the optimization problem logistic regression solves

$$w^{\text{MLE}} = \arg \max_w \sum_i \left\{ y_i w^T x_i - \log \left(1 + e^{w^T x_i} \right) \right\}$$

- How about SVM?

HW1: additional remarks

- P3 (b) Implement the decision tree in Python as described in the book (handles only the binary features) but use the entropy instead of the classification error. Implement an option of `max_depth` so the trained tree will have depth at most `max_depth` (in our case, it corresponds to only considering at most `max_depth` features). Make sure to email your code to csc580homeworks@gmail.com so that I can run it.

Use the data in the book (Table 1) while taking the rating 2/1/0 as positive and -1/-2 as negative. Train your decision tree with your code with `max_depth = 2`. Report your tree along with the following information (in whatever form a person can reasonably comprehend)

- What are the branching questions at each node?
- What are the uncertainty scores at each node?
- Show the predicted label for each leaf node.

- Sanity checks:

- Ensure that `max_depth=2` feature is implemented
- Does your root node's feature have a high information score?

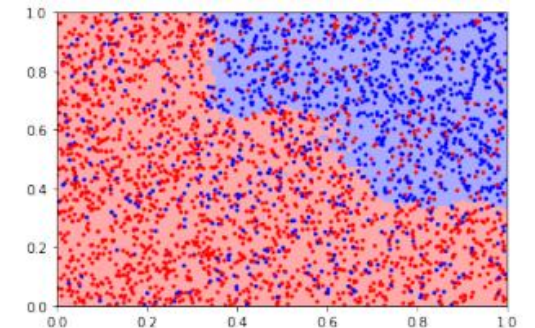
HW1: additional remarks

- P4 Let us define the data distribution \mathcal{D} consisting of two-dimensional features $X \in \mathbb{R}^2$. Each dimension of $X = (X_1, X_2)^T$ is uniformly distributed on the unit interval: $X_1 \sim \text{Uniform}[0, 1]$ and $X_2 \sim \text{Uniform}[0, 1]$. Let $i(X)$ be 1 if $X_1 < 1/3$, 2 if $X_1 \in [1/3, 2/3)$, and 3 if $X_1 \geq 2/3$. Define $j(X)$ similarly for the second dimension X_2 (i.e., replace X_1 above by X_2). Furthermore, the labels are binary with $\mathbb{P}(Y = 1 \mid X = x) = A_{i(x), j(x)}$ and,

$$A = \begin{pmatrix} .1 & .2 & .2 \\ .2 & .4 & .8 \\ .2 & .8 & .9 \end{pmatrix}$$

($A_{i,j} \in \mathbb{R}$ is the entry of matrix A at i -th row, j -th column) Throughout, we abuse notation and use \mathbb{P} for both probability of events and the density function for continuous random variables.

- (a) Bayes optimal classifier
 - Should be an indicator function of a subset of $[0, 1]^2$
 - See solution guide



HW1: additional remarks

(b) Plot the *learning curve* for nearest-neighbor classification. Let $k = 4$. Define $\mathcal{M} = \{10, 30, 100, 300, 1000, 3000\}$. Call the following one ‘trial’:

- Draw 3000 fresh data points from \mathcal{D} and call it S .
- Then, for each $m \in \mathcal{M}$, choose the first m data points from S , train a k -NN classifier with them, and then evaluate its test set error.

Perform 5 trials, compute the average test set error, and report the plot of ‘test error rate’ vs m . In the same plot, plot a horizontal line that shows the Bayes error rate so we know how close we get to the Bayes error rate.

- **Sanity check:**

- The test error of your k -NN classifiers should never be lower than Bayes error rate (why?)

Backup

Basis Functions

- A **basis function** can be any function of the input features X
- Define a set of m basis functions $\phi_1(x), \dots, \phi_m(x)$
- Fit a linear regression model in terms of basis functions,

$$y = \sum_{i=1}^m w_i \phi_i(x) = w^T \phi(x)$$

- Regression model is *linear in the basis transformations*
- Model is *nonlinear in the data X*

Kernel Functions

*A **kernel function** is an inner-product of some basis function computed on two inputs*

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

A consequence is that kernel functions are non-negative real-valued functions over a pair of inputs,

$$\kappa(x, x') \in \mathbb{R} \qquad \kappa(x, x') \geq 0$$

Kernel functions can be interpreted as a measure of distance between two inputs

Kernel Functions

Example The *linear basis* $\phi(x) = x$ produces the kernel,

$$\kappa(x, x') = \phi(x)^T \phi(x') = x^T x'$$

It is often easier to directly specify the kernel rather than the basis function...

Example Gaussian kernel models similarity according to an unnormalized Gaussian distribution,

$$\kappa(x, x') = \exp\left(-\frac{1}{2\sigma^2}(x - x')^2\right)$$

Note Despite the name, this is **not** a Gaussian probability density.

Also called a *radial basis function* (RBF)

Kernel Functions

Given *any* set of data $\{x_i\}_{i=1}^n$ a necessary and sufficient condition of a valid kernel function is that the $n \times n$ **gram matrix**,

$$\mathbf{K} = \begin{pmatrix} \kappa(x_1, x_1) & \kappa(x_1, x_2) & \dots & \kappa(x_1, x_n) \\ \kappa(x_2, x_1) & \kappa(x_2, x_2) & \dots & \kappa(x_2, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ \kappa(x_n, x_1) & \kappa(x_n, x_2) & \dots & \kappa(x_n, x_n) \end{pmatrix}$$

Is a *symmetric positive semidefinite matrix*.

Kernel Ridge Regression

Kernel representation requires inversion of NxN matrix

Primal

$$\Phi = \begin{pmatrix} 1 & \phi_1(x_1) & \dots & \phi_M(x_1) \\ 1 & \phi_1(x_2) & \dots & \phi_M(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \phi_1(x_N) & \dots & \phi_M(x_N) \end{pmatrix}$$

$$w = (\underbrace{\Phi^T \Phi + \lambda I}_{\text{MxM Matrix Inversion}})^{-1} \Phi^T y$$

MxM Matrix Inversion
O(M³)

Dual

$$\mathbf{K} = \begin{pmatrix} \kappa(x_1, x_1) & \kappa(x_1, x_2) & \dots & \kappa(x_1, x_n) \\ \kappa(x_2, x_1) & \kappa(x_2, x_2) & \dots & \kappa(x_2, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ \kappa(x_n, x_1) & \kappa(x_n, x_2) & \dots & \kappa(x_n, x_n) \end{pmatrix}$$

$$y(x) = \mathbf{k}(x)^T (\underbrace{\mathbf{K} + \lambda I}_{\text{NxN Matrix Inversion}})^{-1} \mathbf{y}$$

NxN Matrix Inversion
O(N³)

Number of training data N greater than basis functions M

k -NN classification: pseudocode

- Training is trivial: store the training set
- Test:

Algorithm 3 KNN-PREDICT(\mathbf{D}, K, \hat{x})

list	→	1: $S \leftarrow []$	
		2: for $n = 1$ to N do	
append to list	→	3: $S \leftarrow S \oplus \langle d(x_n, \hat{x}), n \rangle$	// store distance to training example n
		4: end for	
sort in first coordinate	→	5: $S \leftarrow \text{SORT}(S)$	// put lowest-distance objects first
		6: $\hat{y} \leftarrow 0$	
		7: for $k = 1$ to K do	
		8: $\langle \text{dist}, n \rangle \leftarrow S_k$	// n this is the k th closest data point
		9: $\hat{y} \leftarrow \hat{y} + y_n$	// vote according to the label for the n th training point
		10: end for	
Majority vote of $\{y_i: i \in N(x)\}$	→	11: return SIGN(\hat{y})	// return +1 if $\hat{y} > 0$ and -1 if $\hat{y} < 0$

- Time complexity (assuming distance calculation takes $O(d)$ time)
 - $O(m d + m \log m + k) = O(m(d + \log m))$

Inductive Bias

Training



class A



class B

Test



How would you label the test examples?