

ENGR-E 511; ENGR-E 399

Machine Learning for Signal Processing

Module 12:

Adaptive Basis Function Models

Minje Kim

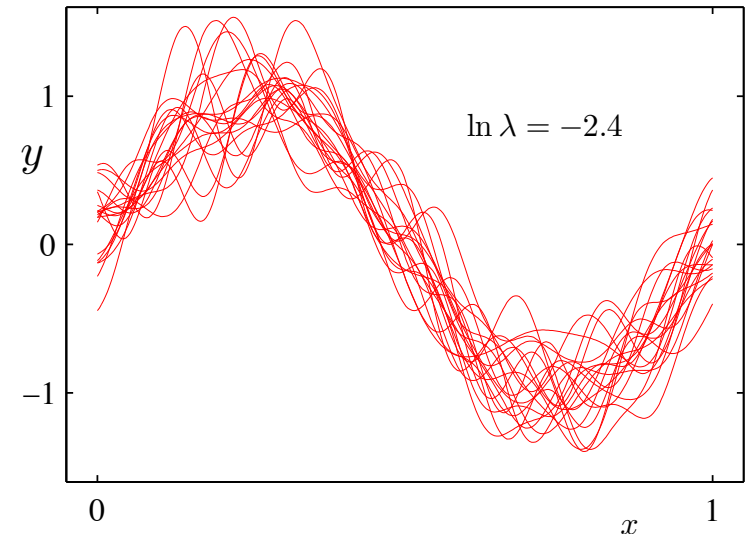
Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>



INDIANA UNIVERSITY
**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**

Adaptive Basis Function Models

- Really, nothing but weighted sum of features

○ Let's go against the argument about the Kernel methods

○ Basis functions

- Feature transformation function $\phi_m(\mathbf{x})$
- In adaptive basis function models we explicitly learn this function from data
 - Instead of using kernels

○ Adaptive basis functions?

- First you assume that there are M such basis functions

$$f(\mathbf{x}) = w_0 + \sum_{m=1}^M w_m \phi_m(\mathbf{x})$$

- The basis function is parameterized and learned from data

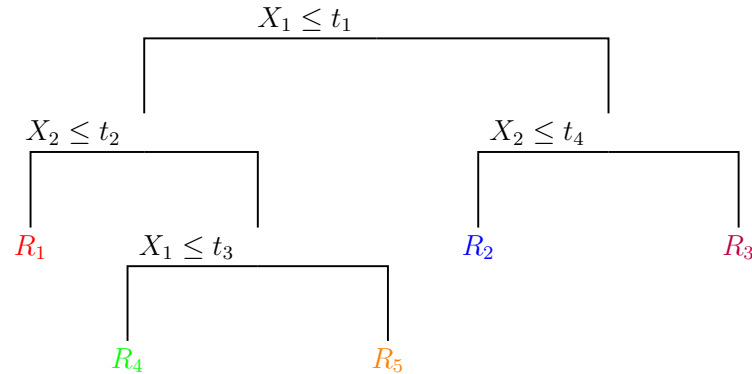
$$\phi_m(\mathbf{x}) = \phi(\mathbf{x}; \mathbf{v}_m)$$

- The entire parameter set

$$\boldsymbol{\theta} = (w_0, w_1, \dots, w_M, \mathbf{v}_1, \dots, \mathbf{v}_M)$$

Decision Trees

- Classification and Regression Trees (CART)



- Axis parallel splits
- Decision trees are an adaptive basis function model

$$f(\mathbf{x}) = \mathbb{E}(y|\mathbf{x}) = \sum_{m=1}^M w_m \mathcal{I}(\mathbf{x} \in \mathcal{R}_m) = \sum_{m=1}^M w_m \phi(\mathbf{x}; \mathbf{v}_m)$$

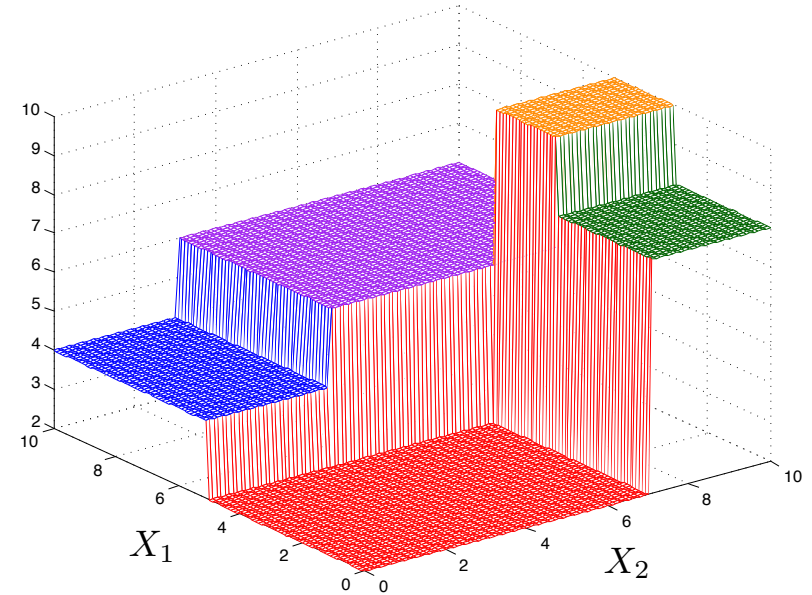
- w_m ?

- Mean response

- \mathbf{v}_m ?

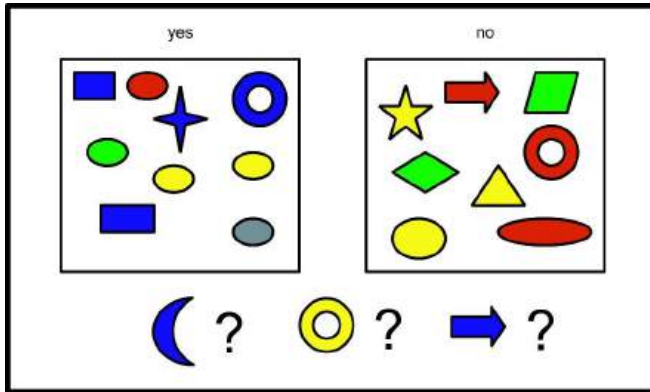
- Which variable to cut, where to cut, and the path from the root to the leaf

$$\phi(\mathbf{x}; \mathbf{v}_m) = \begin{cases} 1 & \text{if } \mathbf{x} \in \mathcal{R}_m \\ 0 & \text{otherwise} \end{cases}$$

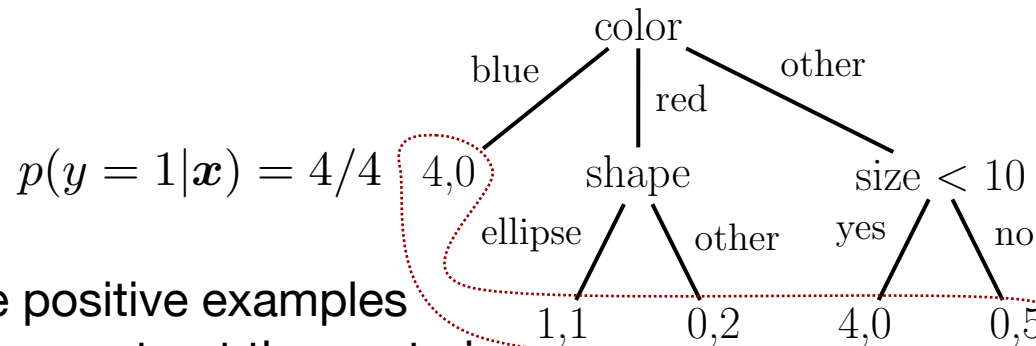


Decision Trees

- CART for classification



D features (attributes)			Label
Color	Shape	Size (cm)	
Blue	Square	10	1
Red	Ellipse	2.4	1
Red	Ellipse	20.7	0



- Fraction of the positive examples can be used to construct the posterior

Decision Trees

- Training

- We need to decide
 - Which feature to threshold
 - Where to threshold
 - Based on the cost minimization
- Optimization

$$(f^*, \tau^*) = \arg \min_{f \in \{1, \dots, F\}} \min_{\tau \in \mathcal{T}_f} \text{cost}(\{\mathbf{x}_t, y_t; x_{ft} \leq \tau\}) + \text{cost}(\{\mathbf{x}_t, y_t; x_{ft} > \tau\})$$

t Sample index

f Features

τ Threshold

\mathcal{T}_f Set of possible threshold

Decision Trees

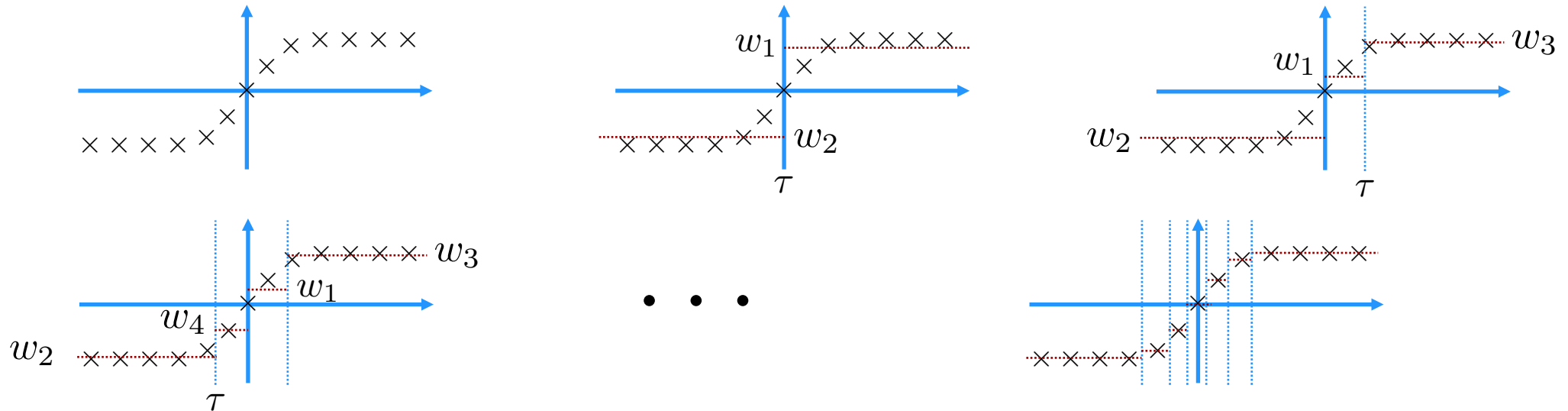
- Regression cost

○ Regression cost $\text{cost}(\mathcal{D}) = \sum_{t \in \mathcal{D}} (y_t - \bar{y})^2$ $\bar{y} = \frac{1}{|\mathcal{D}|} \sum_{t \in \mathcal{D}} y_t$

$t \in \mathcal{D}$ ← Set of examples in the same region (leaf)

○ For 1-d data samples

$$(f^*, \tau^*) = \arg \min_{f \in \{1, \dots, F\}} \min_{\tau \in \mathcal{T}_f} \text{cost}(\{\mathbf{x}_t, y_t; x_{ft} \leq \tau\}) + \text{cost}(\{\mathbf{x}_t, y_t; x_{ft} > \tau\})$$



Decision Trees

- Classification cost

○ Classification cost

$$\hat{\pi}_c = \frac{1}{|\mathcal{D}|} \sum_{t \in \mathcal{D}} \mathbb{I}(y_t = c)$$

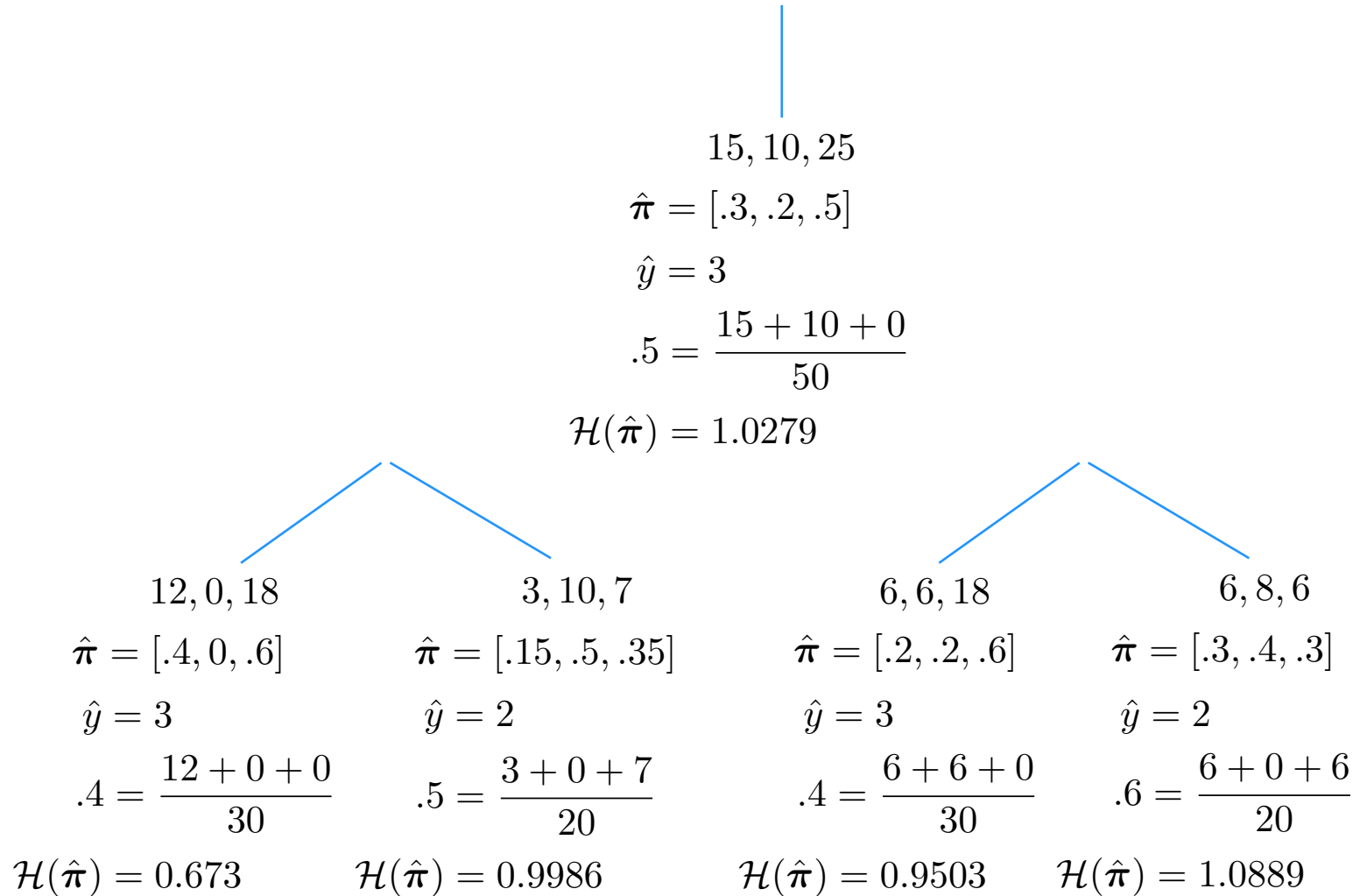
$$\hat{y} = \arg \max_c \hat{\pi}_c$$

□ Misclassification rate

$$\frac{1}{|\mathcal{D}|} \sum_{t \in \mathcal{D}} \mathbb{I}(y_t \neq \hat{y}) = 1 - \hat{\pi}_{\hat{y}}$$

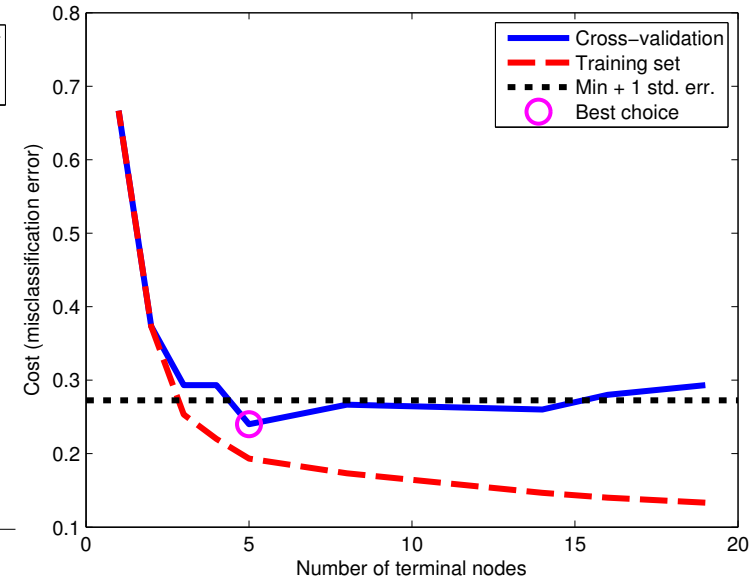
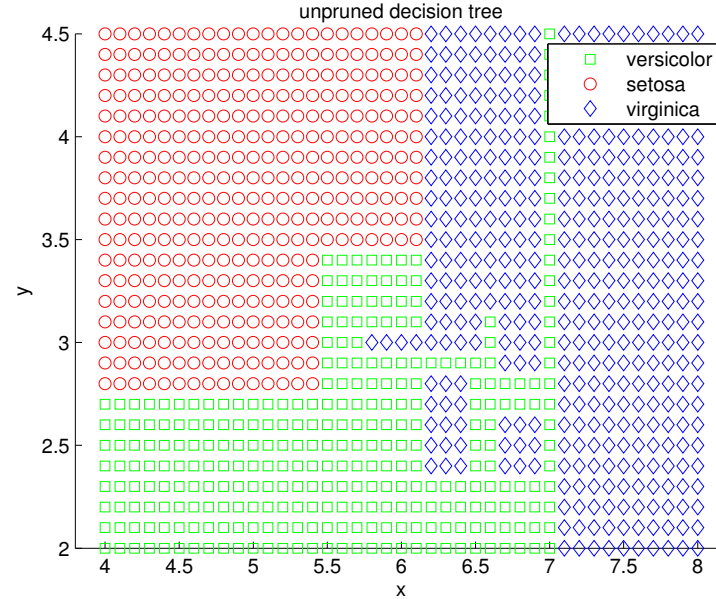
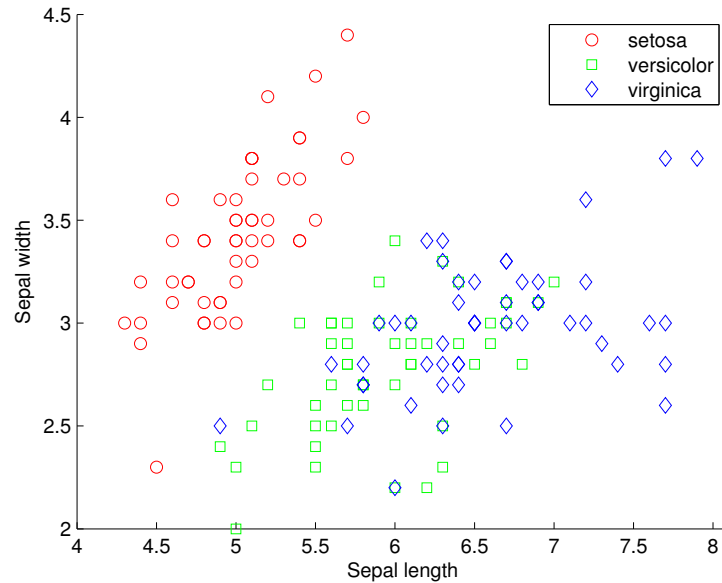
□ Entropy

$$\mathcal{H}(\hat{\pi}) = - \sum_{c=1}^C \hat{\pi}_c \log \hat{\pi}_c$$



Decision Trees

- Iris example

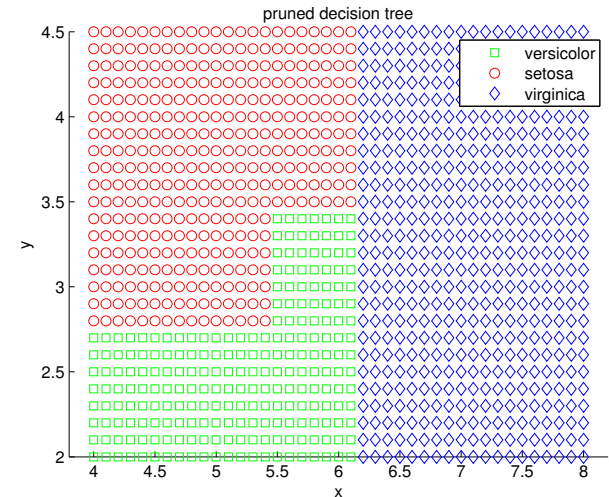
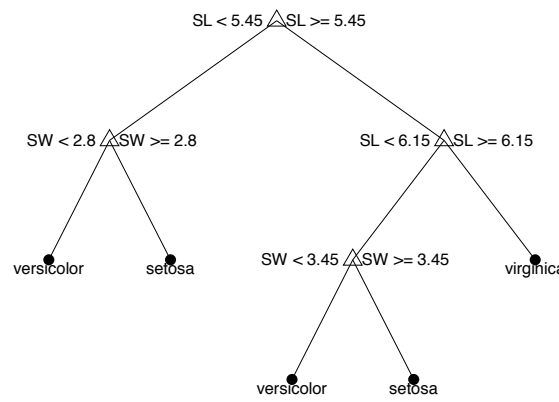
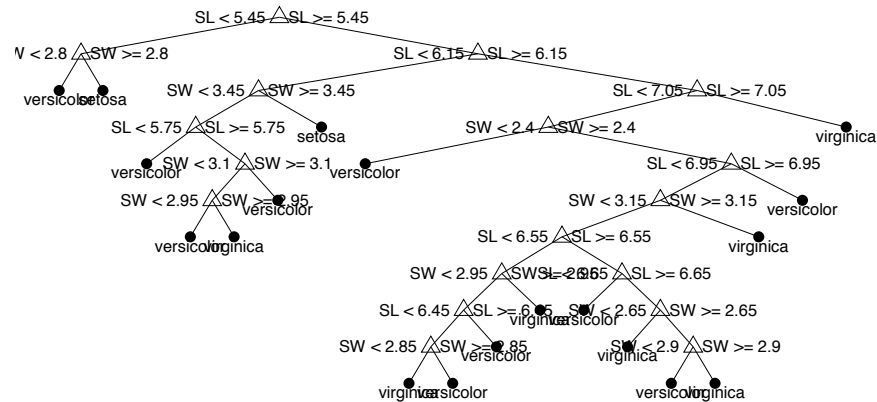
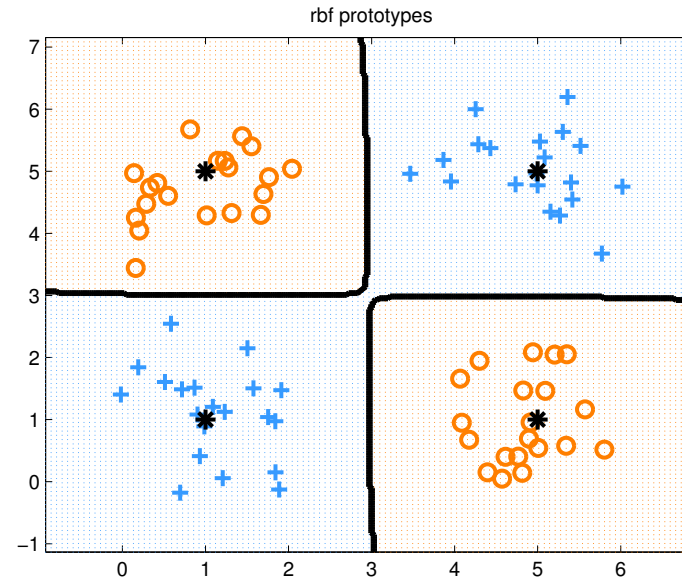


○ Overfitting?

Decision Trees

- Pruning

- Too complex trees can overfit
- To stop growing at some point?
 - e.g. when the cost doesn't go down after a split
 - Doesn't work if there's a confusing case
- Pruning
 - Grow a “full” tree, and then prune



Decision Trees

- Pros and cons of decision trees

○ Pros

- Easy to interpret (e.g. for medical diagnosis)
- Can handle discrete input
- Robust to monotone transformation and scaling (e.g. log)
- Comes with feature selection
- Works well with large datasets
- Easy to handle missing variables

○ Cons

- There are other outperforming models
 - The greedy construction algorithm is not very optimal
- Unstable—small changes in the top node propagate down to the leaf nodes
- High variance



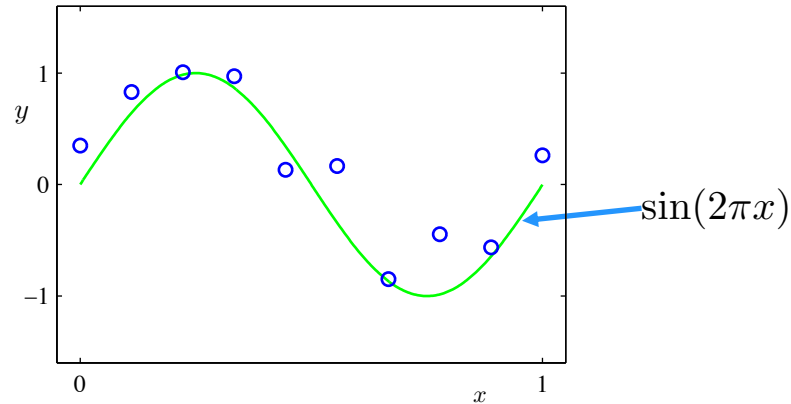
Overfitting and Regularization

- Polynomial curve fitting

- Minimizing an error function

$$\mathcal{E}(\mathbf{w}) = \frac{1}{2} \sum_t (f(x_t, \mathbf{w}) - y_t)^2$$

$$f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \mathbf{w}^\top \mathbf{x}$$



- Therefore,

$$\arg \min_{\mathbf{w}} \mathcal{E}(\mathbf{y} || \mathbf{w}^\top \mathbf{X}) = (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{y} \quad \mathbf{y} = [y_1, y_2, \cdots, y_T]^\top$$

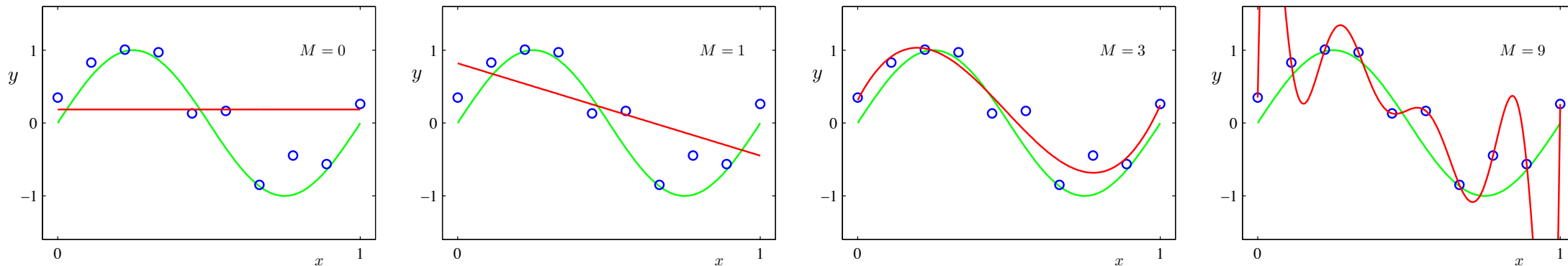
- Question: what's the right order of polynomial, M ?

- The more the better?

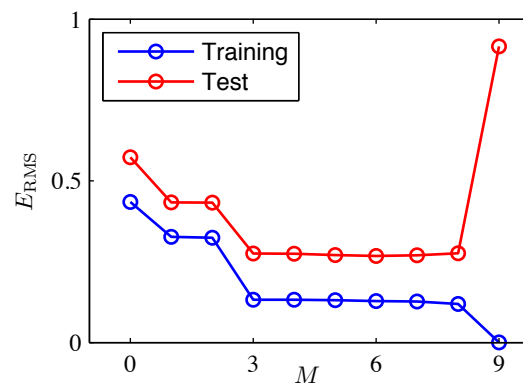
Overfitting and Regularization

- Overfitting

○ The more the better?



- Too complex models tend to **overfit**
 - i.e. The model cares too much about the training error
 - Too complex models don't generalize well
- Generalization?
- A model trained on a dataset (training set) should work well in the other data set (testing set)



Overfitting and Regularization

- Preventing overfitting—early stopping

- Early stopping

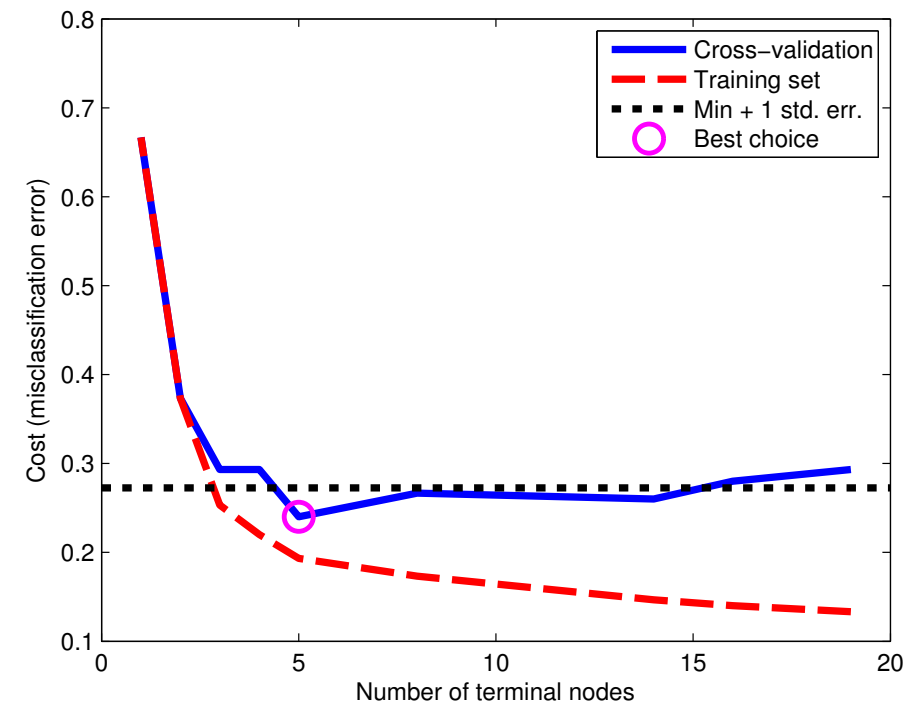
- Check on the simulated test error and stops earlier than the convergence (of the training error)

- N-fold Cross validation: simulate the testing environment using training data

- Divide the training set into N exclusive subsets

	Frames		
Features	1 st fold	2 nd fold	3 rd fold
	Train	Train	Test
	Train	Test	Train
	Test	Train	Train

- N different train-validation pairs
 - Each pair is used to train a classifier and to evaluate it
 - Average the N results
 - The average shows the performance of your choice



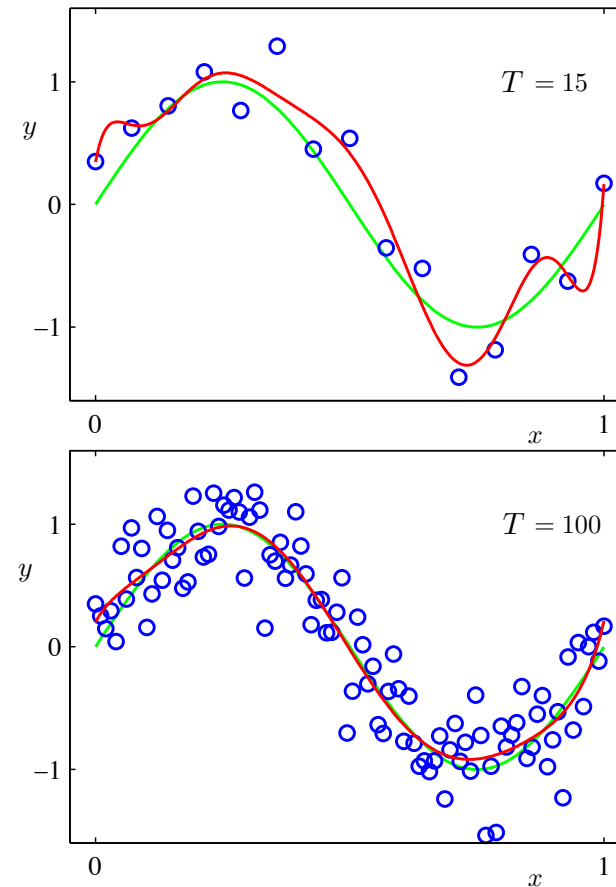
Overfitting and Regularization

- Preventing overfitting—more training samples

○ Weights become larger if the model overfits

○ A big training dataset solves the problem

	$M=0$	$M=1$	$M=3$	$M=9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43



Overfitting and Regularization

- Preventing overfitting—weight decay

- The regularized objective function

$$\mathcal{E}(\mathbf{w}) = \frac{1}{2} \sum_t (f(x_t, \mathbf{w}) - y_t)^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

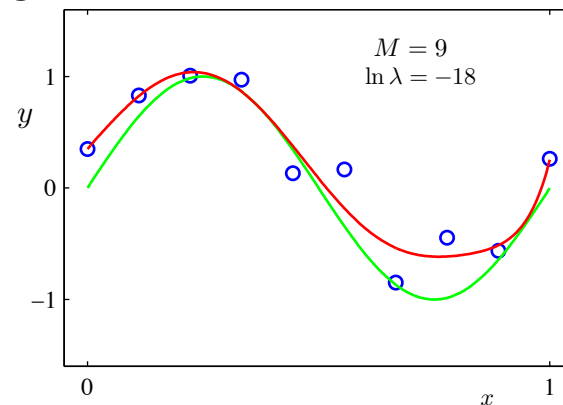
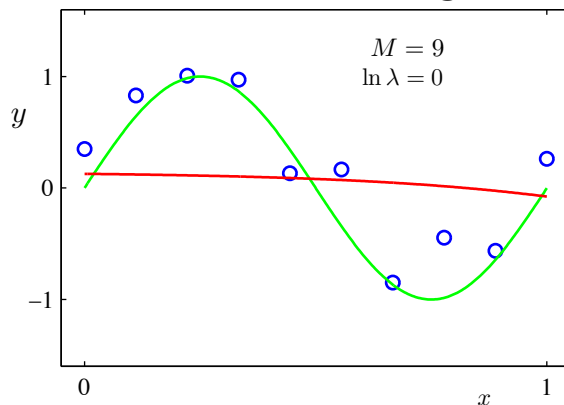
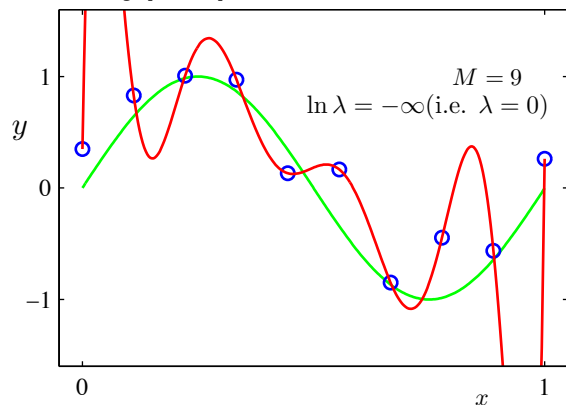
- Keep M large

- You want a complex-enough model to solve your complex problem

- Keep \mathbf{w} small

- You don't want the parameters to be too large

- λ is a hyperparameter that controls the amount of weight decaying



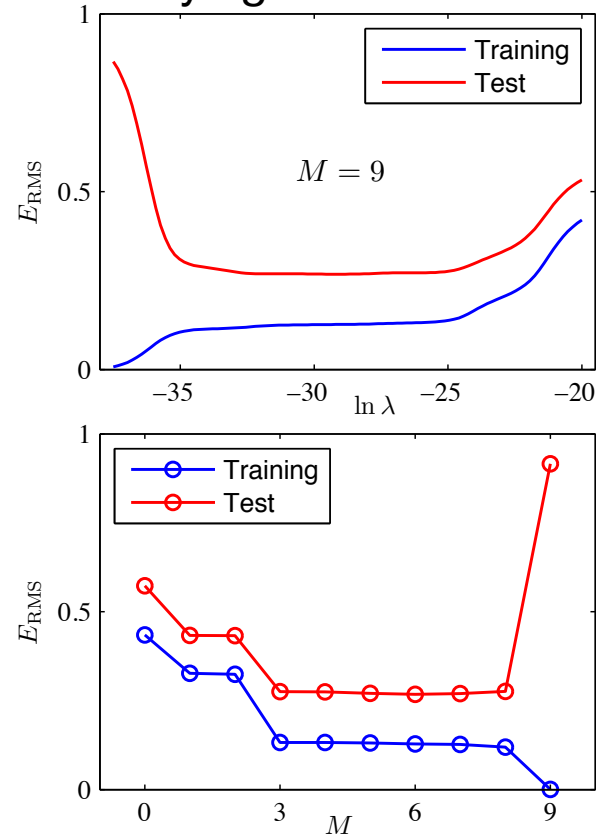
Overfitting and Regularization

- Preventing overfitting—weight decay

- Regularization can decay the weights

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01

- Regularization lets us use a complex model without worrying about overfitting



Overfitting and Regularization

- Preventing overfitting—another takes

- Recall SVM objective function can be seen as a combination of the hinge loss and regularization

$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_t \xi_t \leftrightarrow \sum_t \mathcal{E}(y_t, f(x_t)) + \lambda \|\mathbf{w}\|^2 \quad \mathcal{E}(y_t, f(x_t)) = \xi_t = \begin{cases} 0 & \text{if } y_t f(x_t) \geq 1 \\ 1 - y_t f(x_t) & \text{otherwise} \end{cases}$$

- Bayesian priors can work as a regularizer

- Maximum likelihood

$$P(y|x, \mathbf{w}, \sigma) = \mathcal{N}(y|f(x, \mathbf{w}), \sigma) \leftrightarrow P(\mathbf{y}|\mathbf{x}, \mathbf{w}, \sigma) = \prod_t \mathcal{N}(y_t|f(x_t, \mathbf{w}), \sigma)$$

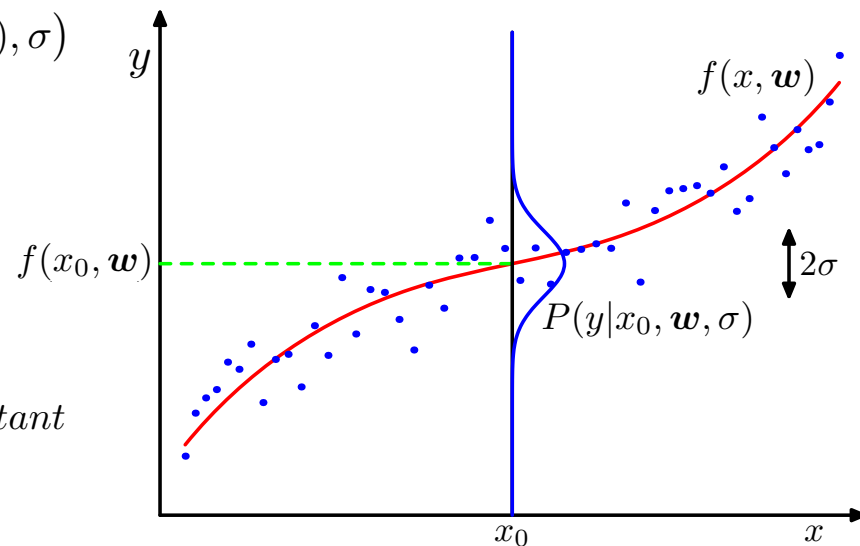
$$\ln P(\mathbf{y}|\mathbf{x}, \mathbf{w}, \sigma) = -\frac{1}{2\sigma^2} \sum_t (f(x_t, \mathbf{w}) - y_t)^2$$

- Prior $P(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma_\pi^2 \mathbf{I}) \propto \exp\left(-\frac{\mathbf{w}^\top \mathbf{w}}{2\sigma_\pi^2}\right)$

- MAP $P(\mathbf{w}|\mathbf{x}, \mathbf{y}, \sigma, \sigma_\pi) \propto P(\mathbf{y}|\mathbf{x}, \mathbf{w}, \sigma)P(\mathbf{w}|\sigma_\pi)$

$$\ln P(\mathbf{w}|\mathbf{x}, \mathbf{y}, \sigma, \sigma_\pi) = -\frac{1}{2\sigma^2} \sum_t (f(x_t, \mathbf{w}) - y_t)^2 - \frac{1}{\sigma_\pi^2} \mathbf{w}^\top \mathbf{w} + \text{constant}$$

- Or to minimize $\frac{1}{2\sigma^2} \sum_t (f(x_t, \mathbf{w}) - y_t)^2 + \frac{1}{\sigma_\pi^2} \mathbf{w}^\top \mathbf{w}$



The Bias-Variance Trade-Off

- Bias-variance decomposition

- Bias-variance decomposition Trained model varies depending on the training set

$$\begin{aligned}
 \mathcal{E}(f(\mathbf{x}; \mathcal{D}) || y) &= (f(\mathbf{x}; \mathcal{D}) - y)^2 \\
 &= (f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] + \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] - y)^2 \\
 &= (f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})])^2 + (\mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] - y)^2 \\
 &\quad + 2(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})])(\mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] - y)
 \end{aligned}$$

Training set
(subsampled from the GT
sample distribution)

Dummy terms

- Expected squared error

$$\mathbb{E}_D \left[(f(\mathbf{x}; \mathcal{D}) - y)^2 \right] = \underbrace{\mathbb{E}_D \left[(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})])^2 \right]}_{\text{Variance:}} + \underbrace{(\mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] - y)^2}_{\text{Bias:}}$$

Expected error;

Assumes multiple models
trained from different subsets
of the original dataset
(by varying \mathcal{D})

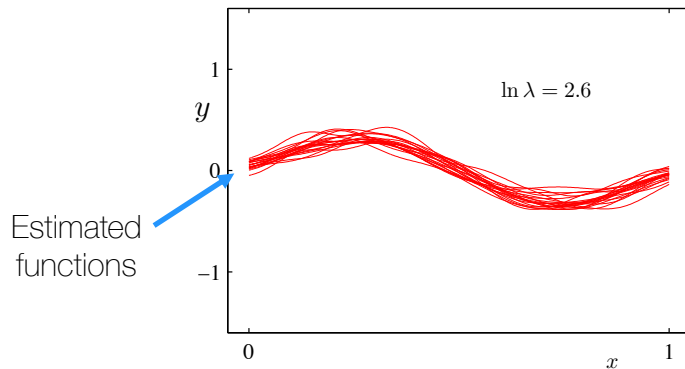
Trained models vary depending on the choice of \mathcal{D} ? How accurate the model is

$$\begin{aligned}
 &+ 2\mathbb{E} \left[(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})])(\mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] - y) \right] \\
 \Leftrightarrow &+ 2(\mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] - y) \mathbb{E} \left[(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})]) \right] \\
 \Leftrightarrow &+ 2(\mathbb{E}_D[f(\mathbf{x}; \mathcal{D})] - y) (\mathbb{E} [f(\mathbf{x}; \mathcal{D})] - \mathbb{E}_D[f(\mathbf{x}; \mathcal{D})]) \xrightarrow{0}
 \end{aligned}$$

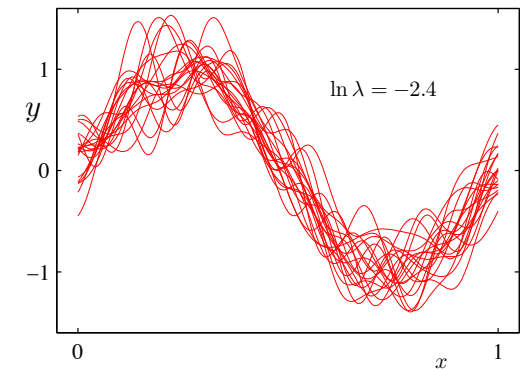
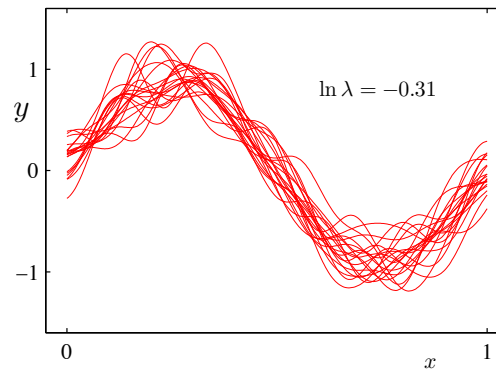
The Bias-Variance Trade-Off

- Averaged model predictions can reduce variance

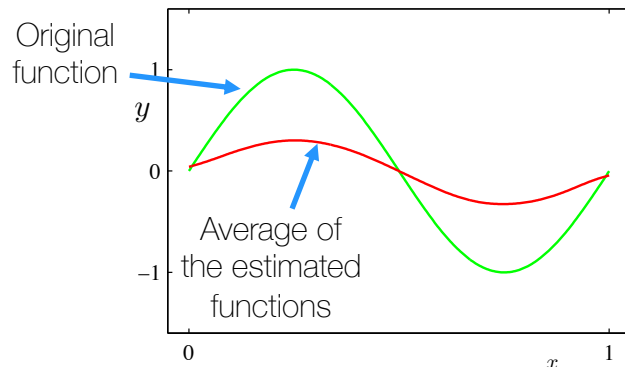
○ 100 models from 100 different \mathcal{D}



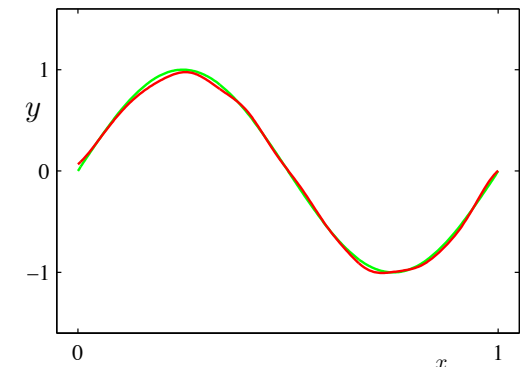
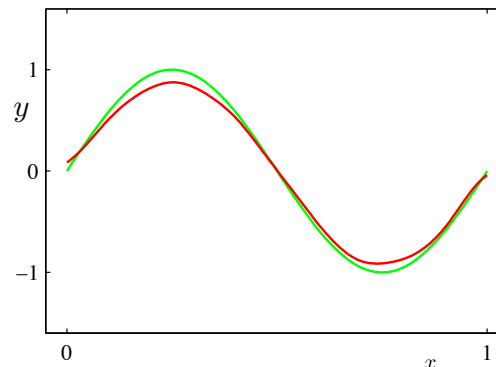
Too much regularization:
low variance, high bias



Too little regularization:
high variance, low bias



Model averaging doesn't help remove bias



Model averaging helps remove variance

The Bias-Variance Trade-Off

- Bootstrap aggregation (or bagging); random forests

- In theory, if you have multiple training datasets,
 - Train multiple complex models and average the results → low variance and low bias
- In practice, you don't have multiple training datasets
- **Bootstrapping**
 - **Subsample from one training dataset with replacement**
- Train m -th model from m -th bootstrap dataset $\sum_{t \in \mathcal{D}_m} \mathcal{E}(\phi_m(\mathbf{x}_t; \mathcal{D}_m) || y_t)$
- The M models construct a committee (**bagging**) $f(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \phi_m(\mathbf{x})$
- **Variance** $\mathbb{E}_{\mathcal{D}} \left[(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x}; \mathcal{D})])^2 \right]$
 - We hope
$$f(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \phi_m(\mathbf{x}) \approx \mathbb{E}_{\mathcal{D}}[f(\mathbf{x}; \mathcal{D})] \quad \text{if } M \rightarrow \infty$$
- **Random forests**: subsample from dataset; subset of variables

AdaBoost (Adaptive Boosting)

- Boosting basics

- Adaptive basis functions are called weak learners in boosting $f(\mathbf{x}) = w_0 + \sum_{m=1}^M w_m \phi_m(\mathbf{x})$

- $\phi_m(\mathbf{x})$ can be a shallow decision tree

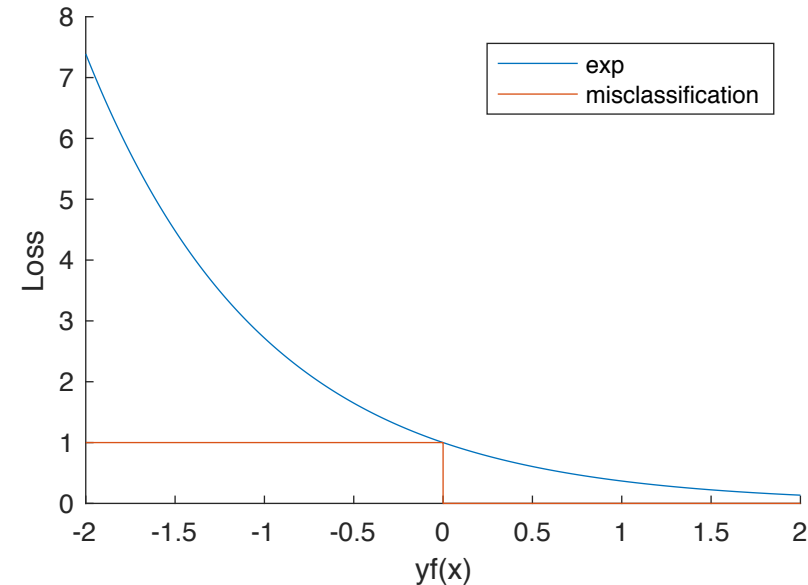
- Or just a perceptron

- The objective function for boosting $\min_f \sum_t \mathcal{E}(y_t, f(\mathbf{x}_t))$

- Exponential loss

$$\mathcal{E}(y, f(\mathbf{x})) = \exp(-yf(\mathbf{x}))$$

- For bipolar binary output $y = \{-1, +1\}$



AdaBoost (Adaptive Boosting)

- Forward stagewise additive modeling

- The exponential loss is magical

- Let's see what happens

- At m -th step of AdaBoost training, we add m -th weak learner to the model

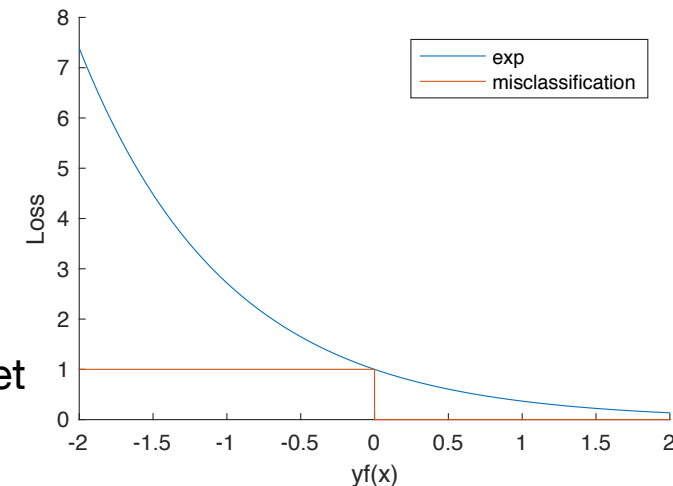
$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m \phi_m(\mathbf{x})$$

- That minimizes the m -th step objective function

$$\begin{aligned}\mathcal{E}_m(\phi) &= \sum_t \exp \left(- y_t (f_{m-1}(\mathbf{x}_t) + \beta_m \phi_m(\mathbf{x}_t)) \right) \\ &= \sum_t \underbrace{\exp \left(- y_t (f_{m-1}(\mathbf{x}_t)) \right)}_{\text{Fixed constants}} \underbrace{\exp \left(- y_t (\beta_m \phi_m(\mathbf{x}_t)) \right)}_{\text{Things to estimate}} \\ &= \sum_t \underbrace{w_{tm}}_{\text{Works like weights over samples}} \exp \left(- \beta_m y_t \phi(\mathbf{x}_t) \right)\end{aligned}$$

- At every step, the weak learner is estimated by using a re-weighted dataset

- Misclassified examples get more weights



AdaBoost (Adaptive Boosting)

- Forward stagewise additive modeling

- Rearrange the objective function

$$\begin{aligned}
 \mathcal{E}_m(\phi) &= \sum_t w_{tm} \exp(-\beta_m y_t \phi(\mathbf{x}_t)) \\
 &= \sum_{t: y_t = \phi(\mathbf{x}_t)} w_{tm} \exp(-\beta_m) + \sum_{t: y_t \neq \phi(\mathbf{x}_t)} w_{tm} \exp(\beta_m) \\
 &= \sum_{t: y_t = \phi(\mathbf{x}_t)} w_{tm} \exp(-\beta_m) + \sum_{t: y_t \neq \phi(\mathbf{x}_t)} w_{tm} \exp(-\beta_m) - \sum_{t: y_t \neq \phi(\mathbf{x}_t)} w_{tm} \exp(-\beta_m) + \sum_{t: y_t \neq \phi(\mathbf{x}_t)} w_{tm} \exp(\beta_m) \\
 &= \sum_{t=1}^T w_{tm} \exp(-\beta_m) + (\exp(\beta_m) - \exp(-\beta_m)) \sum_{t: y_t \neq \phi(\mathbf{x}_t)} w_{tm} \\
 &= \sum_{t=1}^T w_{tm} \exp(-\beta_m) + (\exp(\beta_m) - \exp(-\beta_m)) \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi(\mathbf{x}_t))
 \end{aligned}$$

$\exp(-\beta_m y_t \phi(\mathbf{x}_t)) = \begin{cases} \exp(-\beta) & \text{if } y_t = \phi(\mathbf{x}_t) \\ \exp(\beta) & \text{if } y_t \neq \phi(\mathbf{x}_t) \end{cases}$

- Estimating the m-th weak learner is eventually $\phi_m = \arg \min_{\phi} \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi(\mathbf{x}_t))$

□ For example $\text{cost}(\mathcal{D}) = \sum_{t \in \mathcal{D}} w_{tm} (y_t - \bar{y})^2$

But, you give different weights to the different samples, based on the previous round,

(m-1)-th step

This is a misclassification loss. So, basically you're looking for a simple, but good classifier



AdaBoost (Adaptive Boosting)

- Forward stagewise additive modeling

○ Solve for β

$$\begin{aligned}\mathcal{E}_m(\phi) &= \sum_{t=1}^T w_{tm} \exp(-\beta_m) + (\exp(\beta_m) - \exp(-\beta_m)) \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t)) \\ 0 &= - \sum_{t=1}^T w_{tm} \exp(-\beta_m) + (\exp(\beta_m) + \exp(-\beta_m)) \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t)) \\ 0 &= - \sum_{t=1}^T w_{tm} \exp(-\beta_m) \exp(\beta_m) + (\exp(\beta_m) \exp(\beta_m) + \exp(-\beta_m) \exp(\beta_m)) \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t)) \\ 0 &= - \sum_{t=1}^T w_{tm} + \exp(2\beta_m) \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t)) + \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t)) \\ 0 &= \exp(2\beta_m) \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t)) - \sum_{t=1}^T w_{tm} \mathcal{I}(y_t = \phi_m(\mathbf{x}_t)) \\ \exp(2\beta_m) \sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t)) &= \sum_{t=1}^T w_{tm} \mathcal{I}(y_t = \phi_m(\mathbf{x}_t))\end{aligned}$$

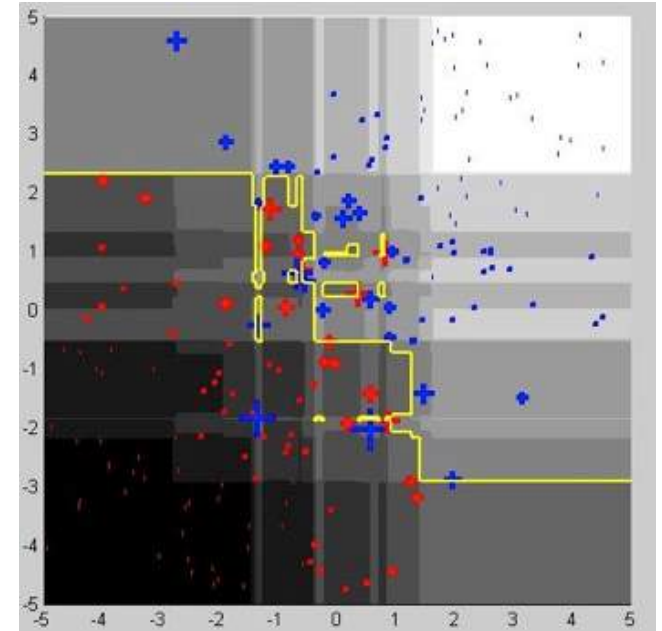
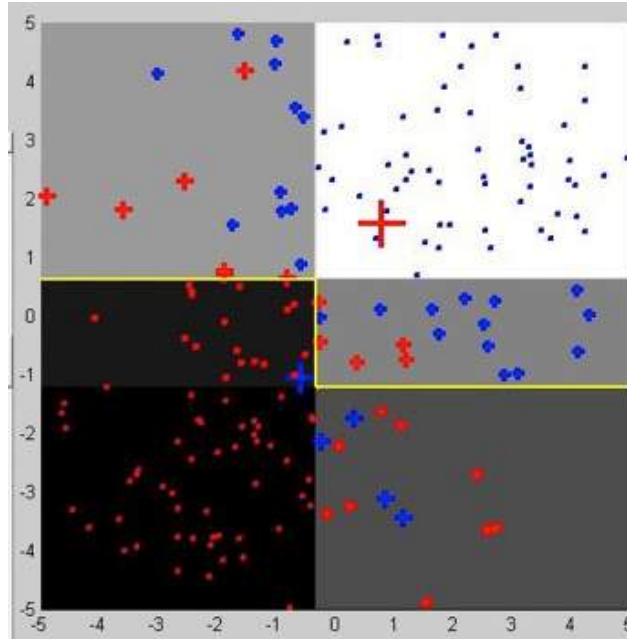
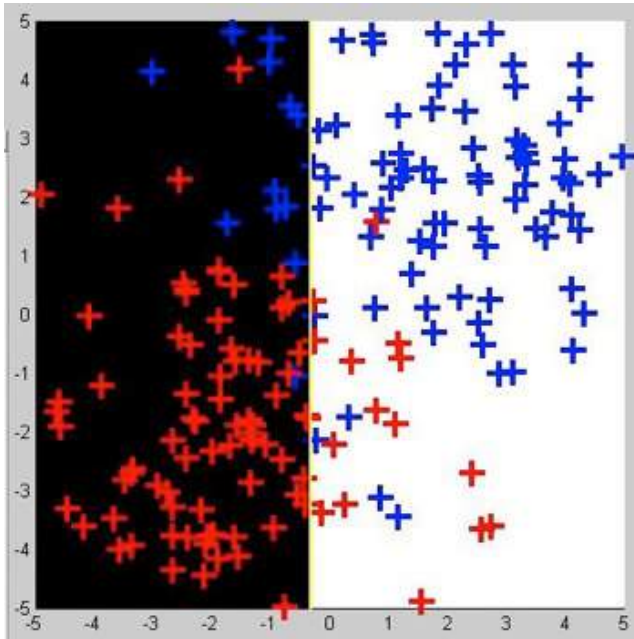
Differentiation

$$\therefore \beta_m = \frac{1}{2} \ln \frac{\sum_{t=1}^T w_{tm} \mathcal{I}(y_t = \phi_m(\mathbf{x}_t))}{\sum_{t=1}^T w_{tm} \mathcal{I}(y_t \neq \phi_m(\mathbf{x}_t))}$$

AdaBoost (Adaptive Boosting)

- Forward stagewise additive modeling

- Now that all the estimation jobs are done, $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m \phi_m(\mathbf{x})$
- Calculation of the new weights $w_{t,m+1} = w_{t,m} \exp(-\beta_m y_t \phi_m(\mathbf{x}_t))$ See slide 22 if you want to prove
- Repeat this process by adding more basis functions

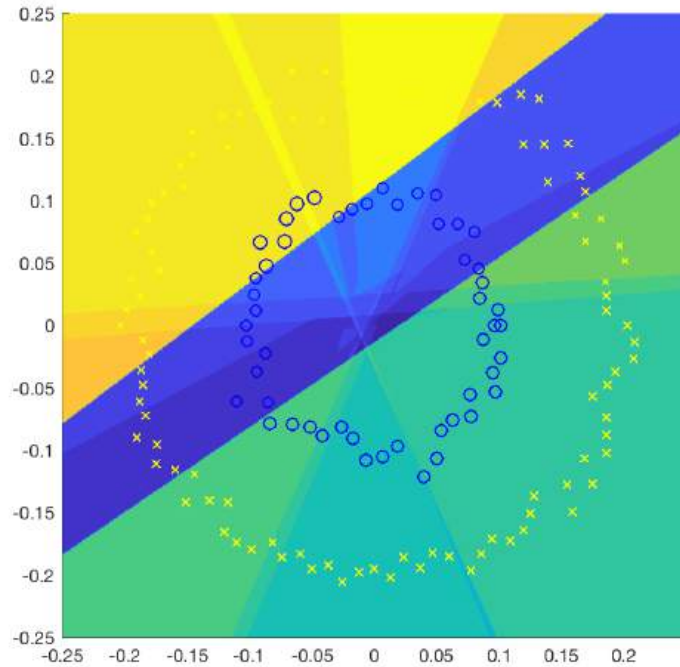


Size of the markers represents their weights

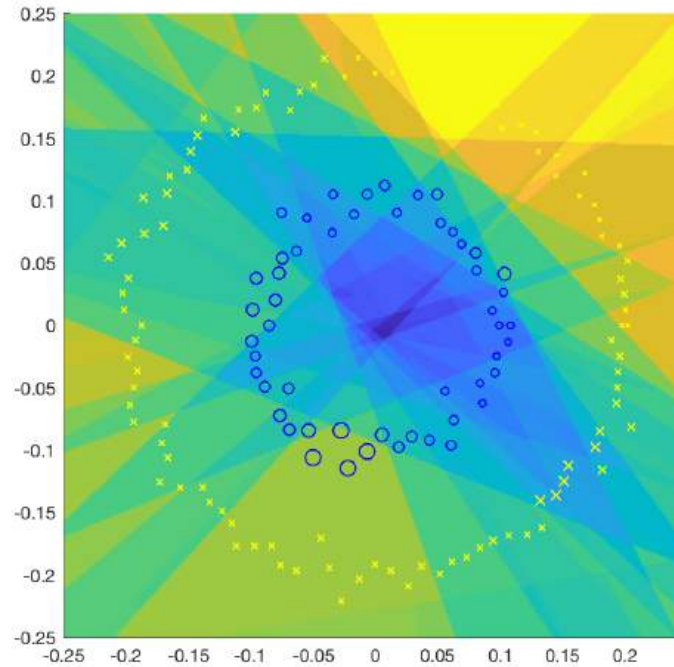
AdaBoost (Adaptive Boosting)

- Concentric

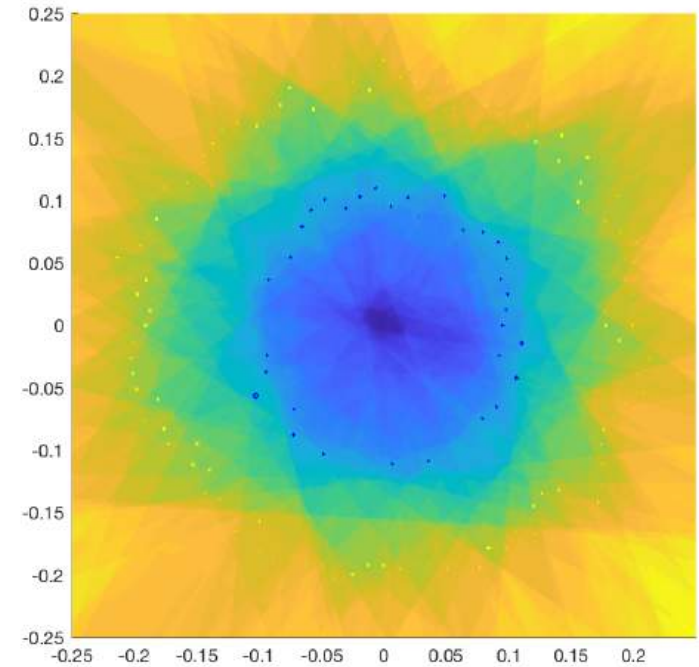
- AdaBoost on the concentric dataset



$m = 20$



$m = 50$



$m = 500$

Deep Neural Networks

- A neural network is an adaptive basis function model

$$f(\mathbf{x}) = w_0 + \sum_{m=1}^M w_m \phi_m(\mathbf{x})$$

One of the Hidden unit outputs in the final layer

Weight associated with the unit

- Deep neural networks does the feature transformation using all the layers up to the last one
 - So, the basis functions can be seen as features
 - Note that boosting sees the basis functions in this way, too
- DNNs tend to give you an unbiased solution
 - In other words, they overfit if the training set is small
 - Needs strong regularization with not enough data
 - Deep learning: big models for big data
 - The popularity of CNN: models are smaller by introducing a more complex operation, convolution
- ENGR E533 “Deep Learning Systems” in Spring 2019

Reading

- Kevin Murphy, “Machine Learning: a Probabilistic Perspective”
 - Chapter 6.4, 16.1-16.4
- Christopher Bishop, “Pattern Recognition and Machine Learning”
 - Chapter 1.1, 3.2





Thank You!

