

CS 340: Machine Learning

Lecture 3: Introduction to Supervised Learning

AD

January 2011

Supervised learning as function fitting

- We are given some training data

$$\mathcal{D} = \{(\mathbf{x}^i, y^i)\}_{i=1}^N$$

- Consider a restricted set of mappings/parametric functions f in *hypothesis class* \mathcal{H}

$$f \in \mathcal{H} : \mathcal{X} \times \Theta \longrightarrow \mathcal{Y},$$

we will predict using

$$\hat{y}(\mathbf{x}) = f(\mathbf{x}; \theta)$$

where $\theta \in \Theta$.

- **Learning:** Given \mathcal{H} , learn parameters θ given \mathcal{D} so that predictions on non-labeled inputs (i.e. test set, real-world data) are as accurate as possible.

- **Empirical/Training error:**

$$\text{Err_Train} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{y}(\mathbf{x}^i) \neq y^i).$$

- **Test error:**

$$\text{Err_Test} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \mathbb{I}(\hat{y}(\mathbf{x}_{\text{test}}^i) \neq y_{\text{test}}^i).$$

- **Generalization error:**

$$\text{Proba}[\text{error}] = \mathbb{E}[\mathbb{I}(\hat{y}(\mathbf{x}) \neq y)] = \int \mathbb{I}(\hat{y}(\mathbf{x}) \neq y) p(\mathbf{x}, y) d\mathbf{x}dy$$

where the expectation is w.r.t the UNKNOWN probability density function $p(\mathbf{x}, y)$ of (\mathbf{x}, y) .

- Empirical/Training error are approximations of generalization error thanks to the Law of Large Numbers.

Law of Large Numbers

- Assume that some \mathcal{Z} -valued random variables $\{z^i\}$ are independent and identically distributed according to a probability density function $p(z)$ then for any function $g: \mathcal{Z} \rightarrow \mathbb{R}$ the law of large numbers states that

$$\lim_{N \rightarrow \infty} \underbrace{\frac{1}{N} \sum_{i=1}^N g(z^i)}_{\text{empirical average}} = \underbrace{\mathbb{E}[g(z)]}_{\text{statistical average}}$$

where the expectation $\mathbb{E}[g(Z)]$ is defined by

$$\mathbb{E}[g(z)] = \int g(z) p(z) dz.$$

- This fundamental theorem states that the limit of the empirical average converges to the statistical average.
- In the supervised learning case, we have simply $z = (\mathbf{x}, y)$ and $g(z) = g(\mathbf{x}, y) = \mathbb{I}(\hat{y}(\mathbf{x}) \neq y)$.

Binary classification: Credit card scoring

- Say you have training data of the following form

Income	Savings	Risk
100	50	Hi
100	100	Lo
50	75	Hi
500	93	Lo

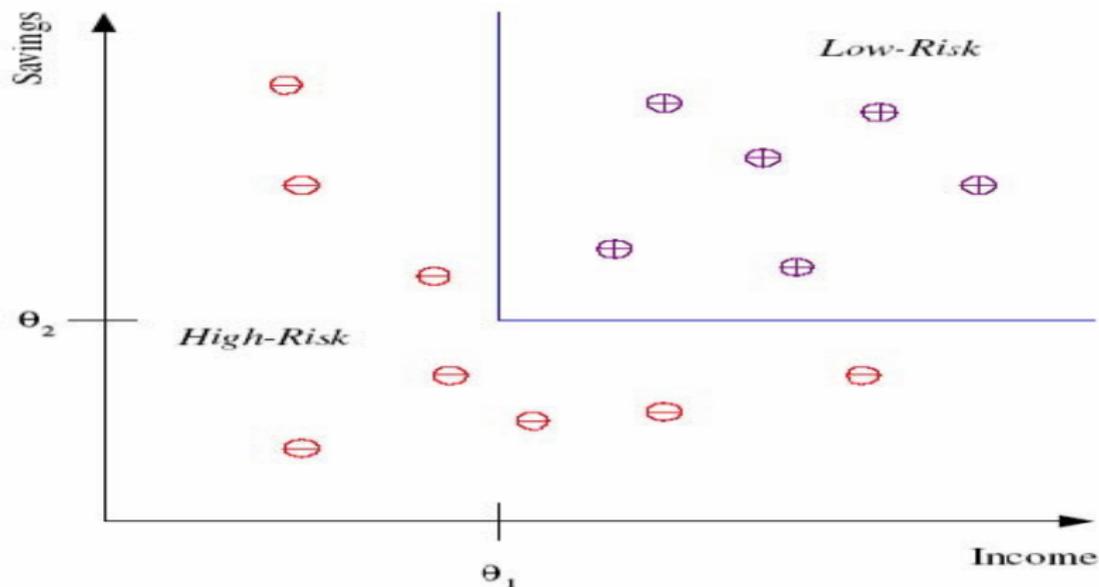
- Test data are of the form

Income	Savings	Risk
98	49	?
100	102	?
400	20	?

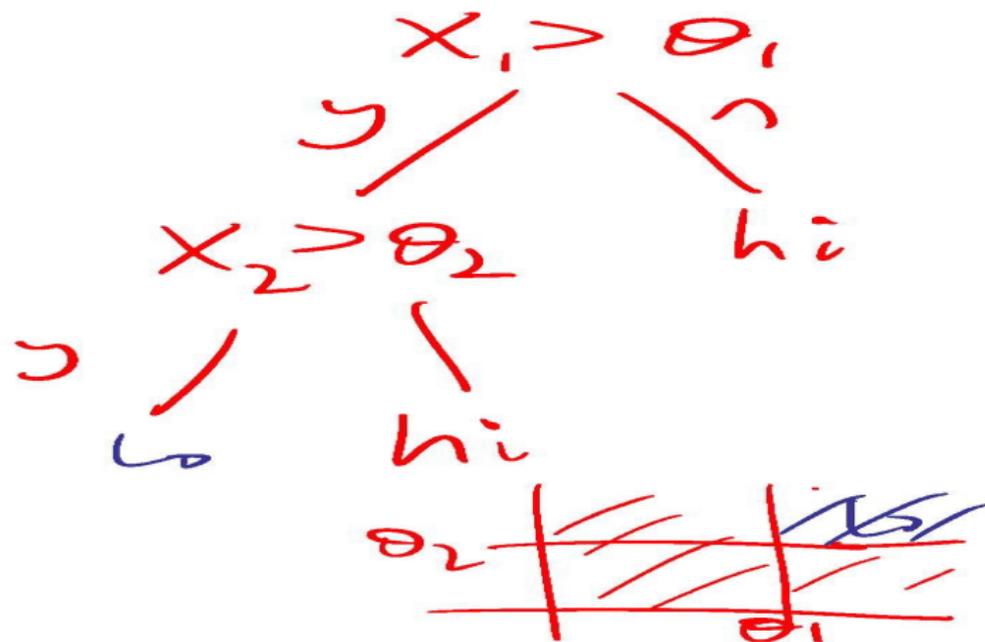
- In this case $\mathbf{x} = (x_1, x_2) \in (\mathcal{X} = \mathbb{R}^2)$ and $\mathcal{Y} = \{\text{high, low}\}$.

Example function

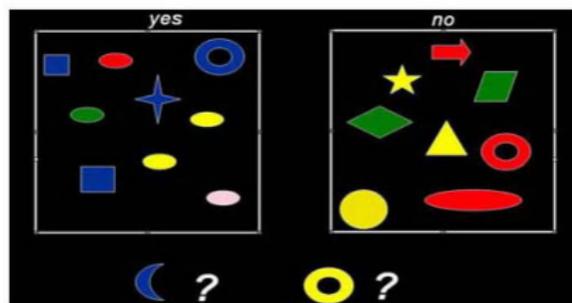
- $f(\mathbf{x}; \theta) = f((\text{income}, \text{savings}); (\theta_1, \theta_2)) = \text{IF } (x_1 = \text{income}) > \theta_1$
AND $(x_2 = \text{savings}) > \theta_2$ THEN low-risk ELSE high-risk.



Decision Trees



Classifying Shapes



p features (attributes)

Training set:

X: n by p
y: n by 1

n cases

Color	Shape	Size
Blue	Square	Small
Red	Ellipse	Small
Red	Ellipse	Large

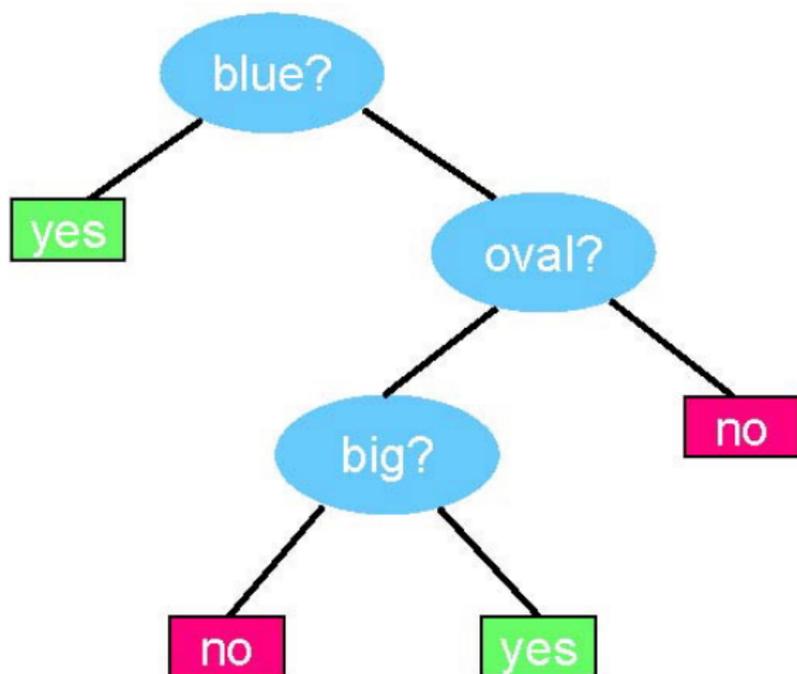
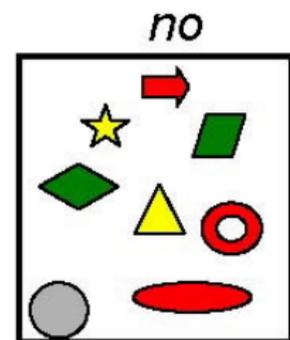
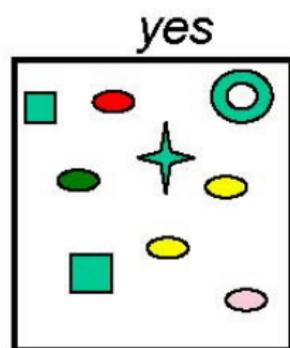
Label
Yes
Yes
No

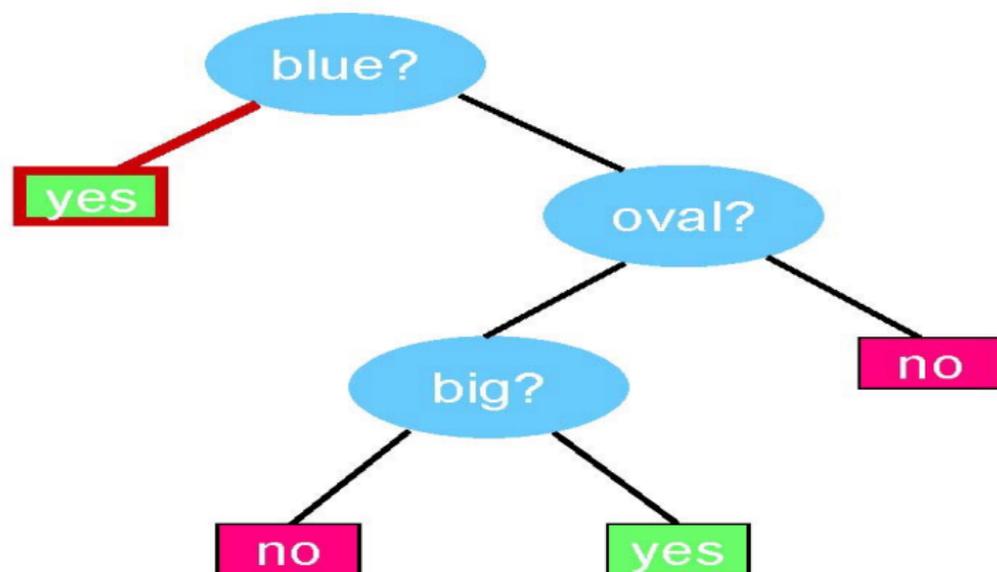
Test set

Blue	Crescent	Small
Yellow	Ring	Small

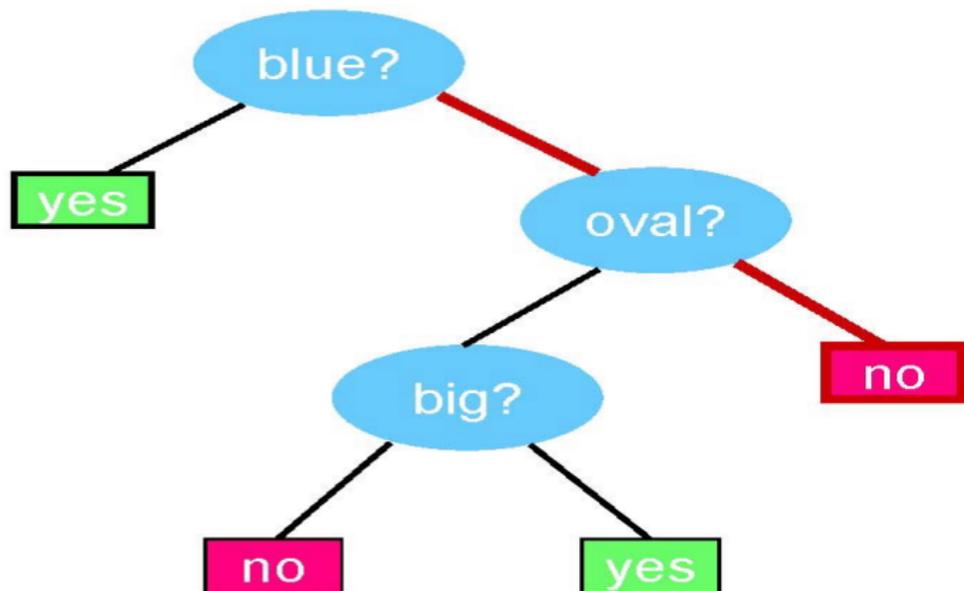
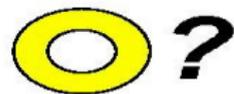
?
?

Hypothesis (Decision Trees)





- Generalization beyond training set can be difficult: Blue \rightarrow yes: Reasonable?



- Generalization beyond training set can be difficult: Yellow small ring → no: Reasonable?

- Decision trees are easy to use and easy to understand, and they can easily handle mixed discrete and continuous inputs.
- They can be very difficult to build (NP hard even in some simple cases).
- Optimal tree structure is sensitive to the particular patterns in the training data: not very stable to perturbations.
- Adding one more example of a certain kind can change the whole tree topology and the model does not generalize very well to novel inputs.
- No longer very popular in machine learning.

What's the right hypothesis class?

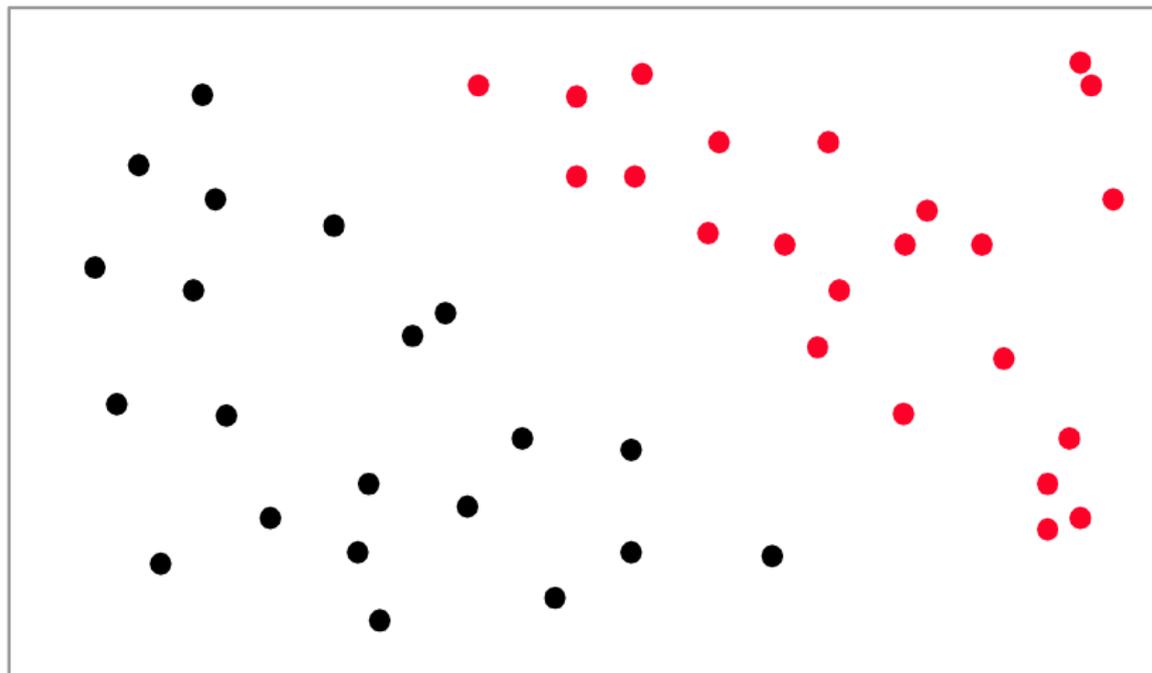


Figure: Training data where $\mathbf{x}^i \in \mathbb{R}^2$ and $y^i \in \{-1, 1\}$

Linearly separable data

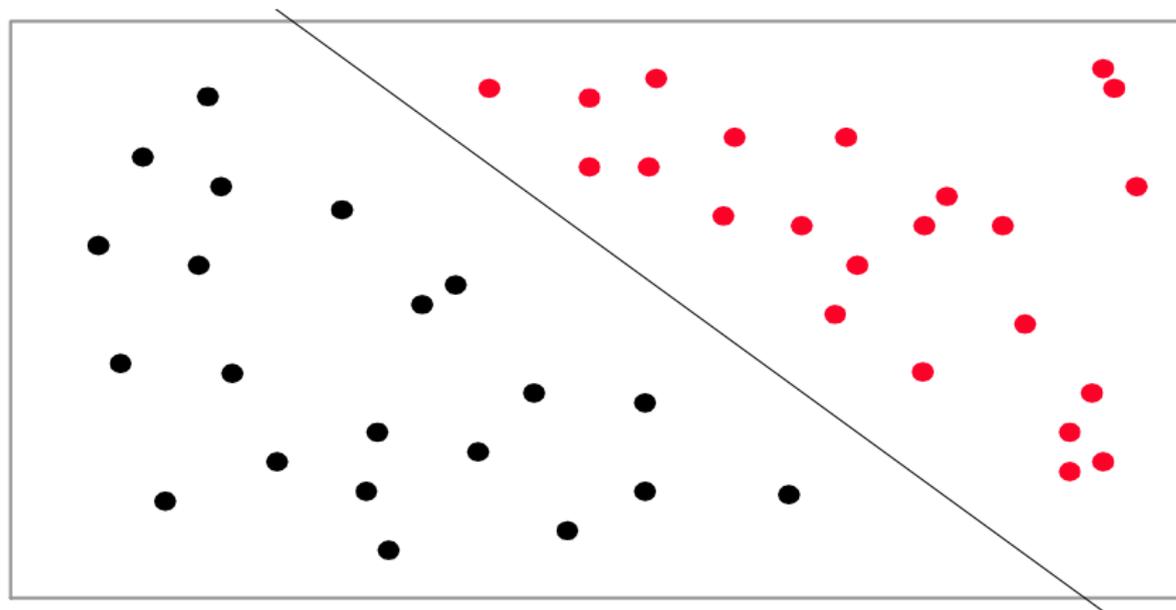
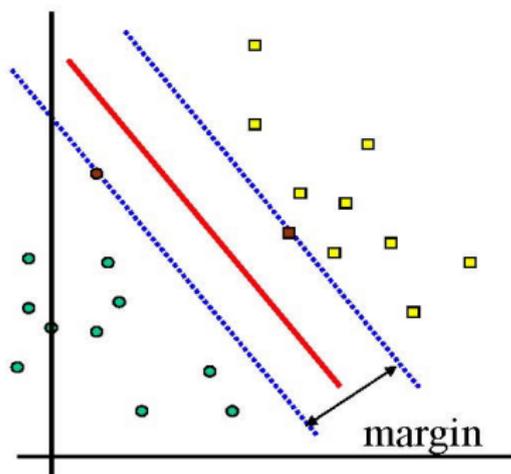
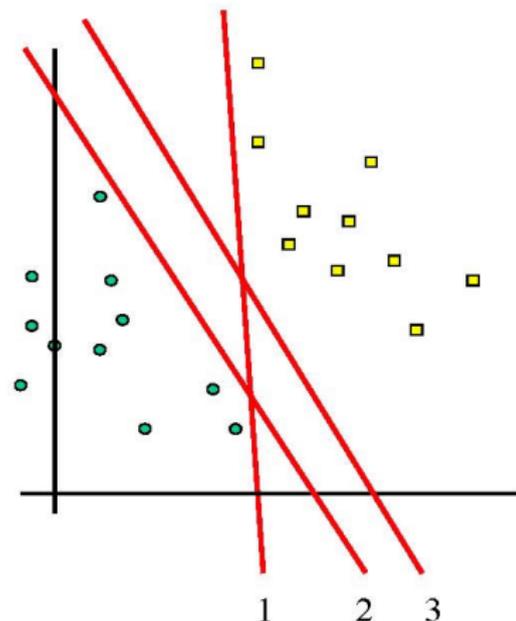


Figure: Linear classifier where $\mathbf{x}^i \in \mathbb{R}^2$ and $y^i \in \{-1, 1\}$

- Linearly separable means if f is a “linear” function of \mathbf{x} , we can perfectly fit the training data.

Which linear hypothesis is better?



- Even if the chosen hypothesis class is properly chosen, we will need to specify a criterion to select θ !

Not linearly separable data

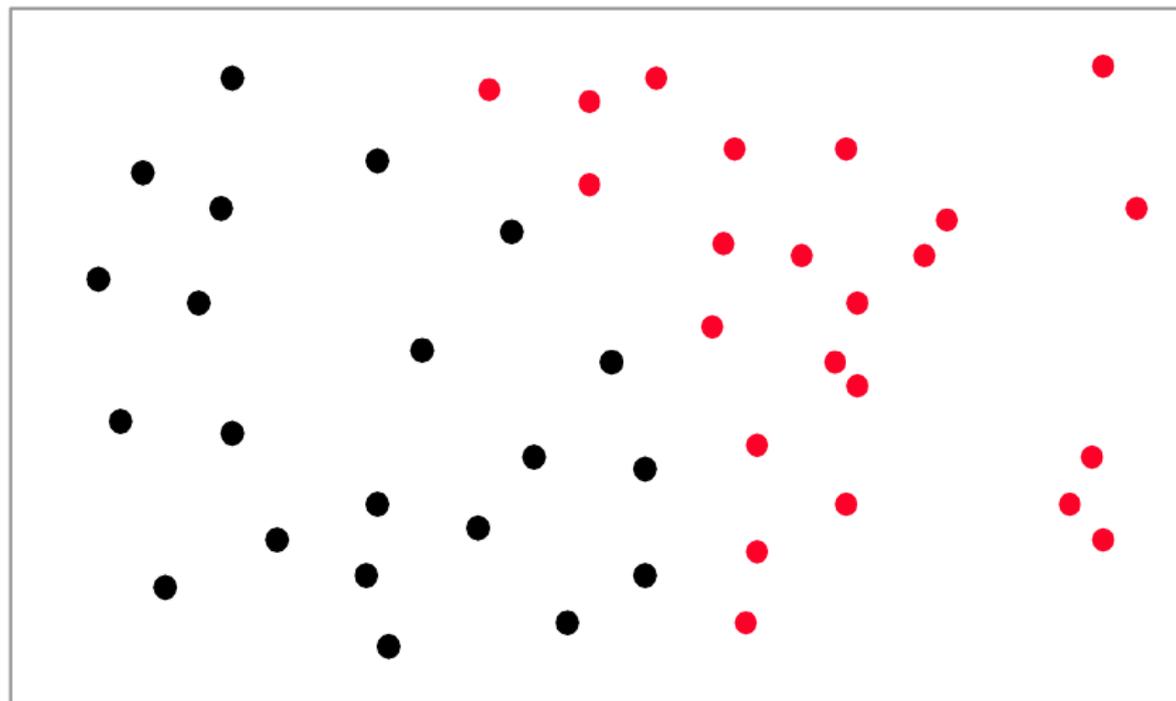
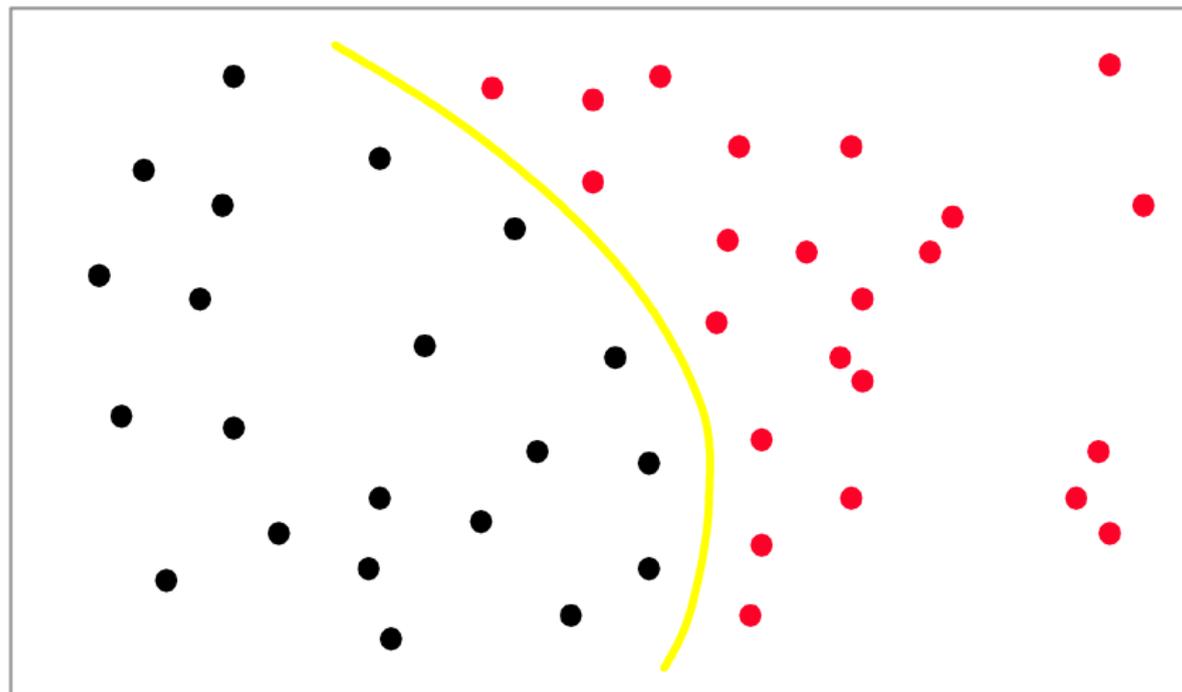


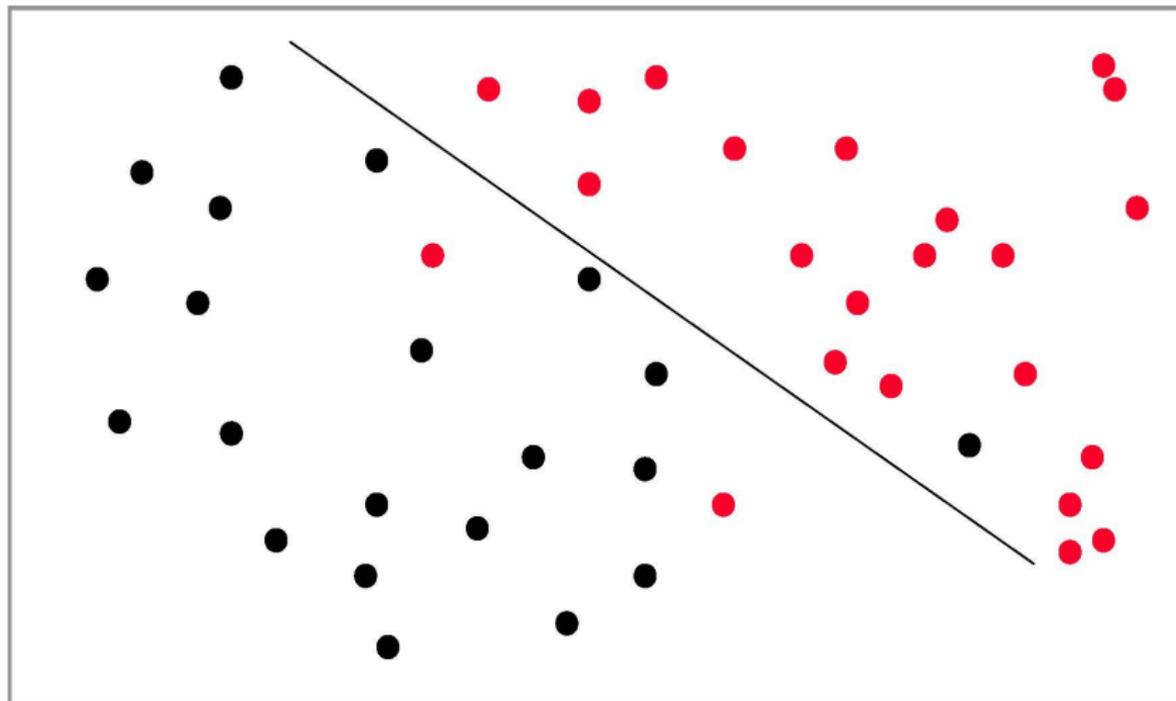
Figure: Non-linearly separable training data where $\mathbf{x}^i \in \mathbb{R}^2$ and $y^i \in \{-1, 1\}$

Quadratically separable

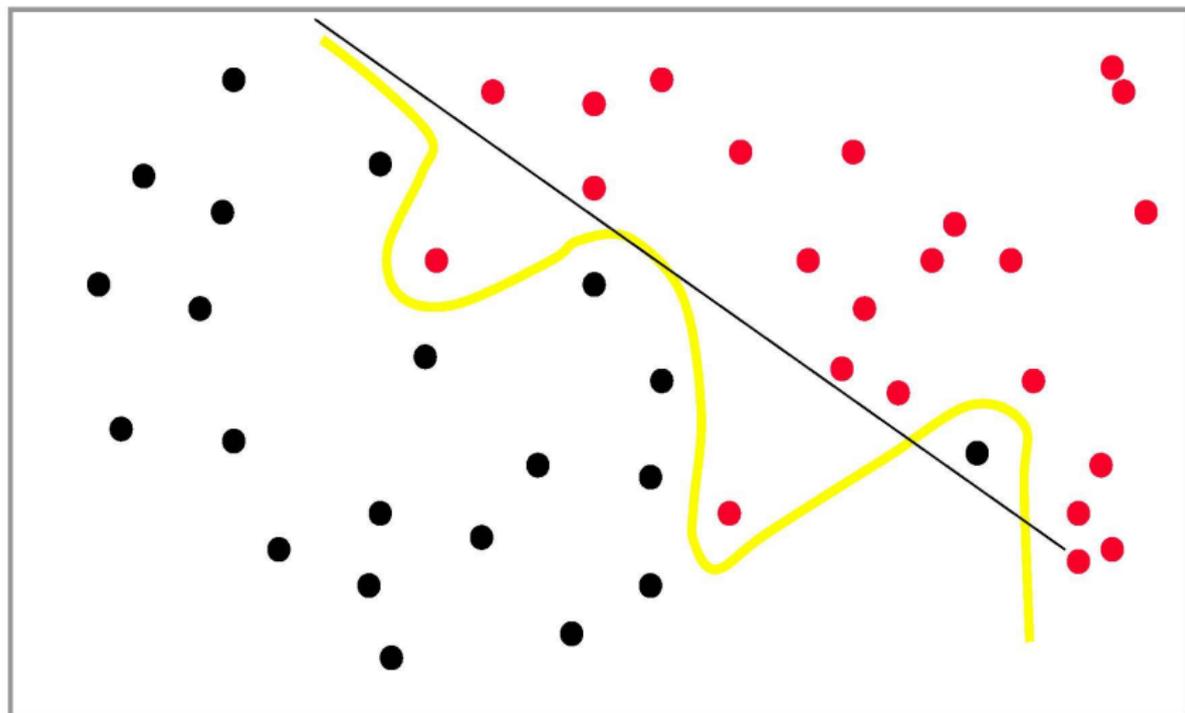


- $f(\mathbf{x}; \theta) = \text{sgn}(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$.

Noisy/Mislabeled data and/or Features non-informative

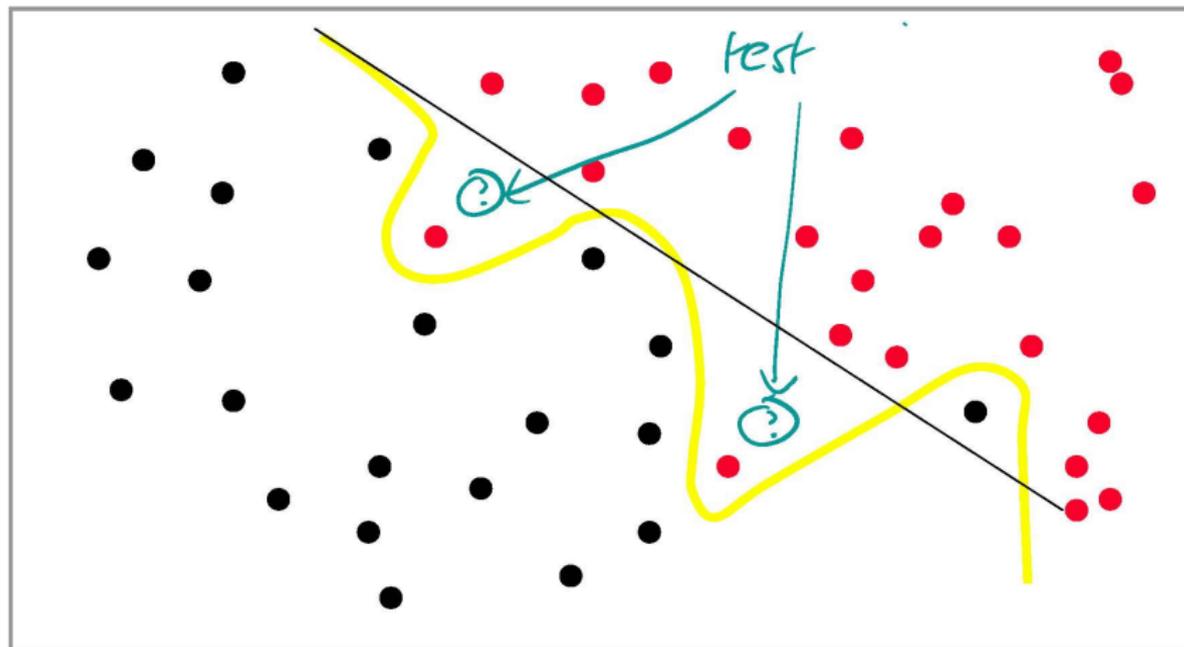


Overfitting

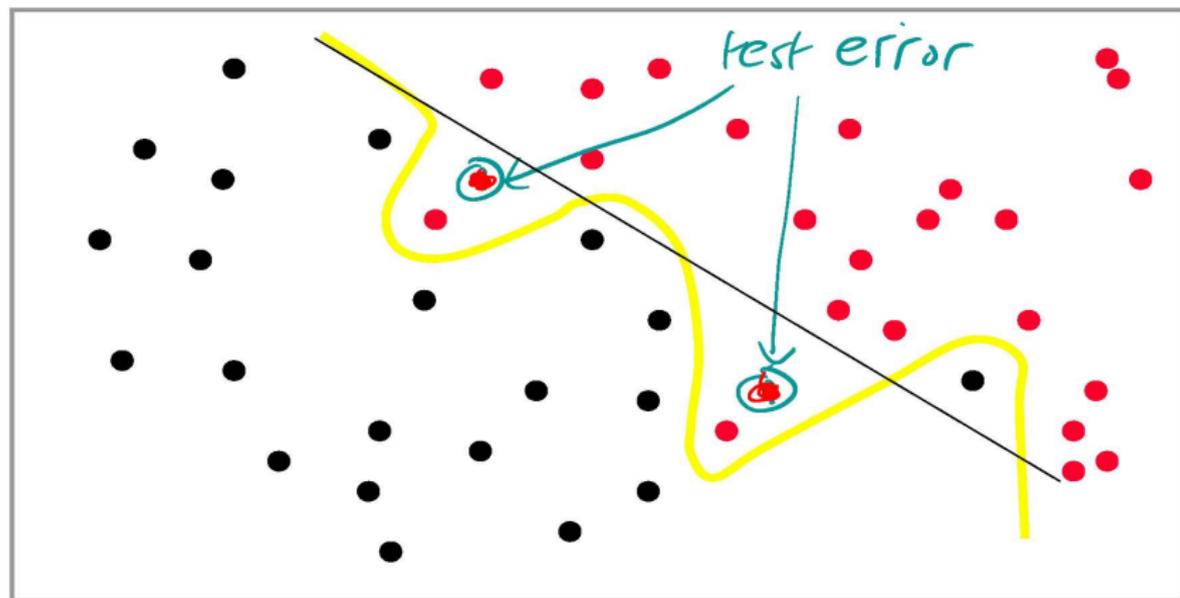


- An overly flexible function might model irrelevant details of training set.

Overfitted functions do not predict well

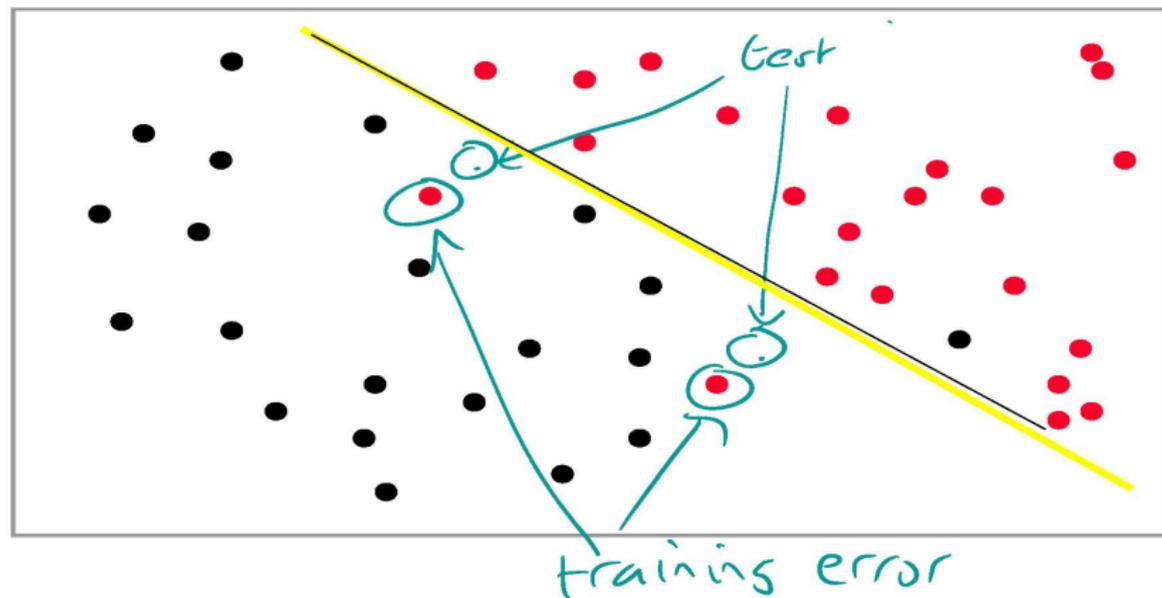


Overfitted functions do not predict well

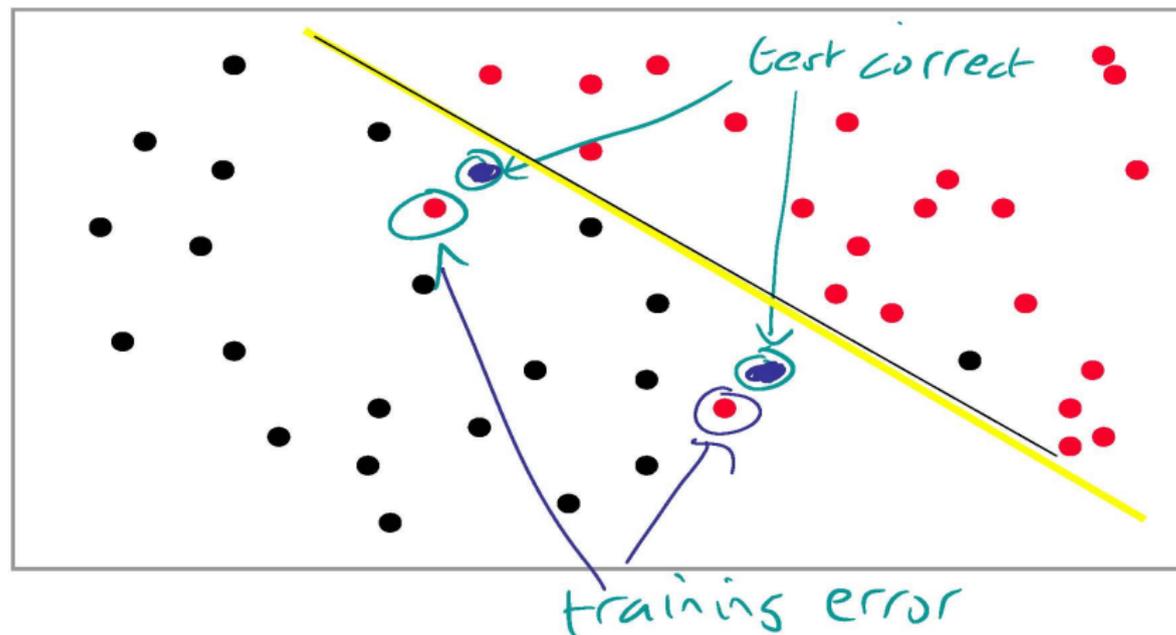


Test points are mis-predicted

Tradeoff simplicity for model fit

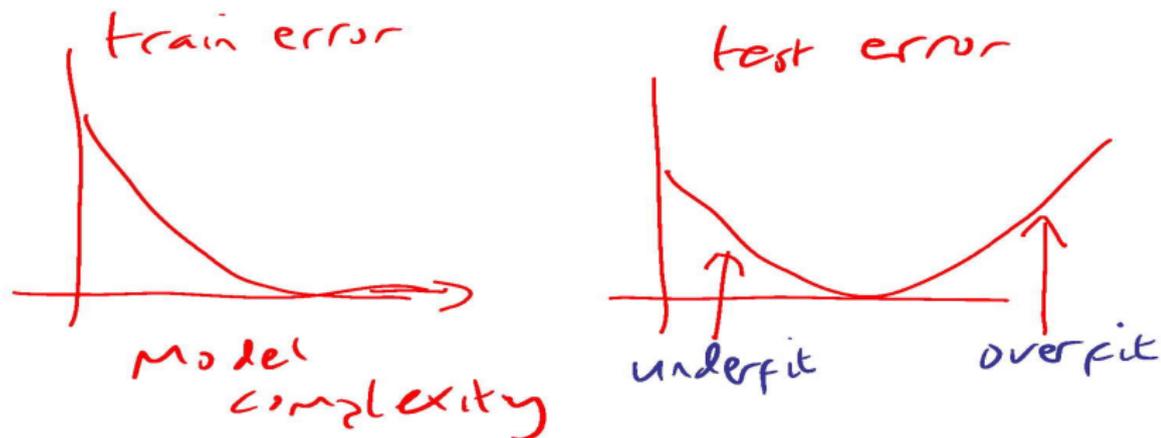


Tradeoff simplicity for model fit



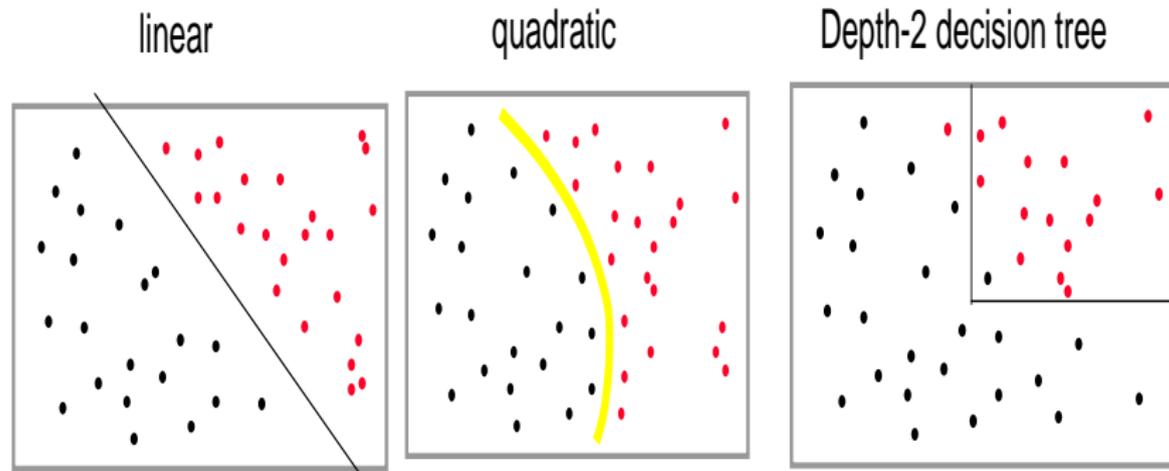
Ockham's razor

- If two models fit the data equally well, pick the simpler one
- In general, since our goal is to predict the test data, we may choose to incur errors on the training set if it results in a simpler function.



Function fitting

- Choose the right hypothesis class \mathcal{H} given training data \mathcal{D} .



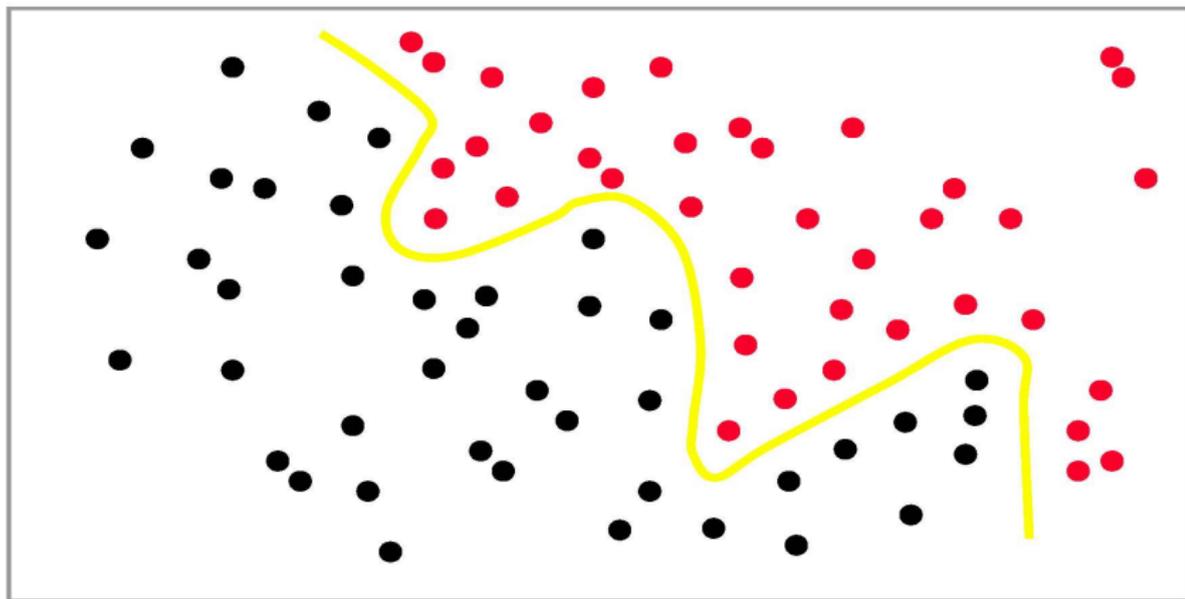
- Learn parameters θ of function $f(\mathbf{x}; \theta)$ given \mathcal{H} and \mathcal{D} ; e.g.

$$f(\mathbf{x}; \theta) = \text{sgn}(\theta^T (1 \ \mathbf{x})) = \text{sgn}(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

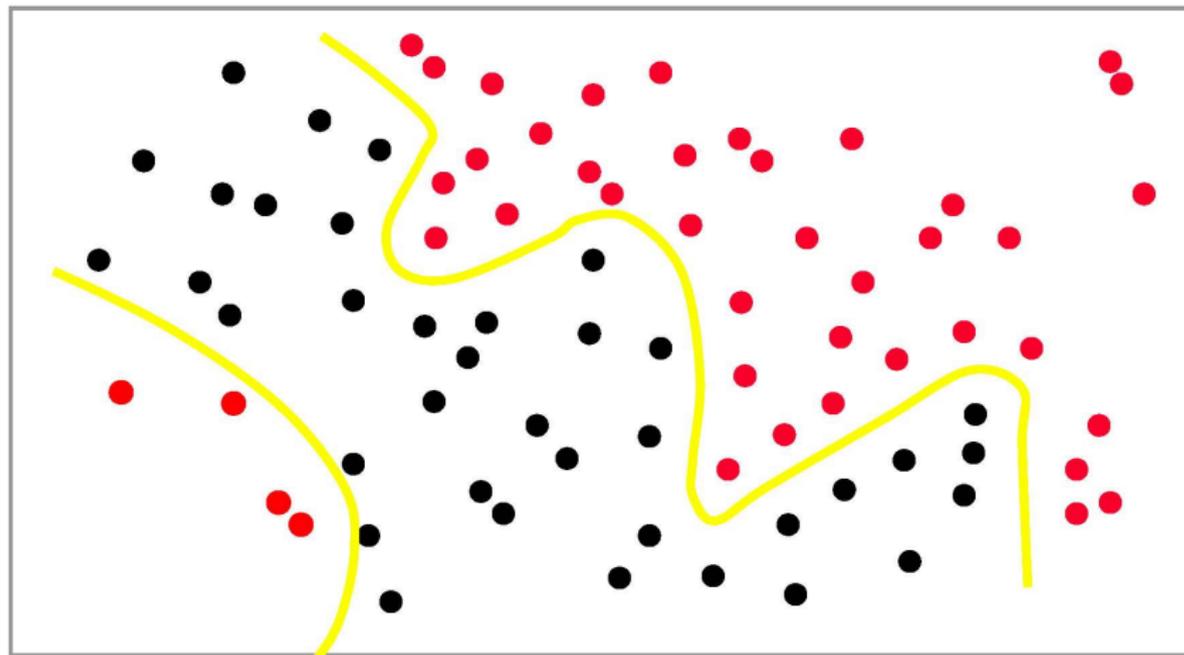
according to a given (typically statistical) criterion.

Hypothesis class depends on data

- More complex function is ok if we have more data, because we have more evidence for it.



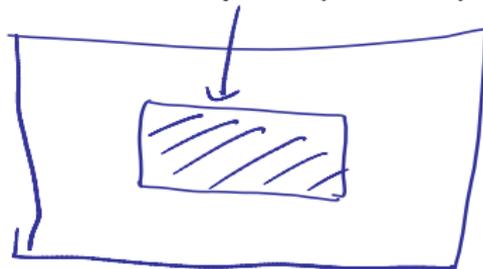
Decision regions might be discontinuous



Another Simple Example

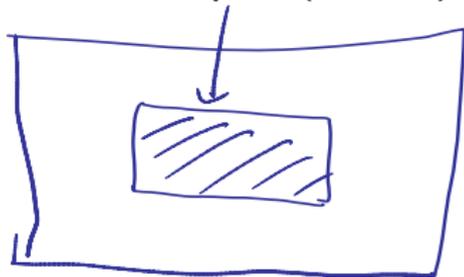
- Learn concept of “healthy levels” of cholesterol x_1 and insulin x_2 from positive and negative examples.
- \mathcal{H} = space of rectangles in \mathbb{R}^2 .
- Assume for the time being that the “true” mapping belongs also to \mathcal{H} .

True concept C (hidden)

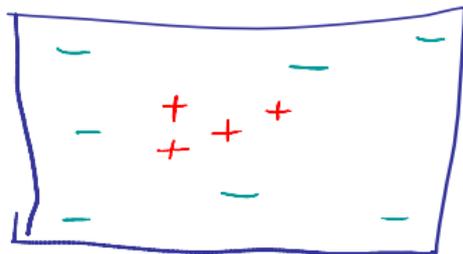


Training Data

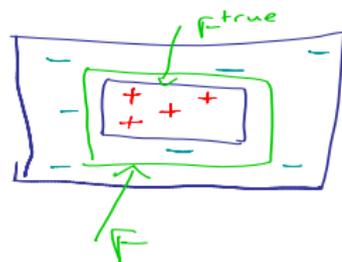
True concept C (hidden)



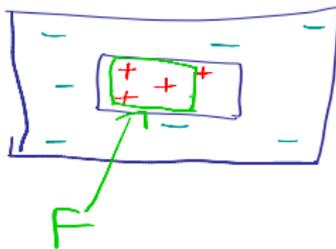
Training data D sampled from C



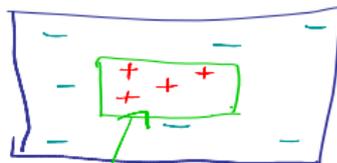
Learning: Inferring hidden concept/function



Too big

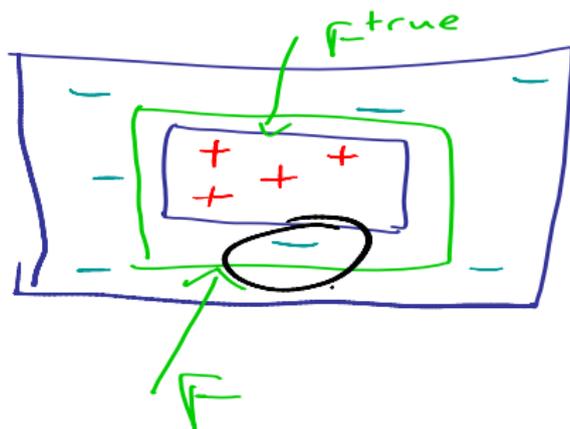


Too small



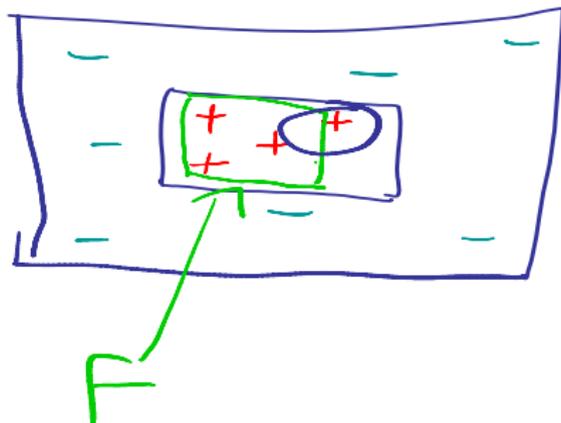
Just right

False Positive



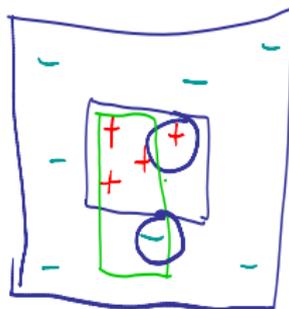
$$N_{\text{FalsePos}} = \sum_{i=1}^N \mathbb{I}(\hat{y}(\mathbf{x}_i) = 1 \wedge y_i = 0)$$

False Negative



$$N_{\text{FalseNeg}} = \sum_{i=1}^N \mathbb{I}(\hat{y}(\mathbf{x}_i) = 0 \wedge y_i = 1)$$

Generalization Error and Empirical Risk Minimization

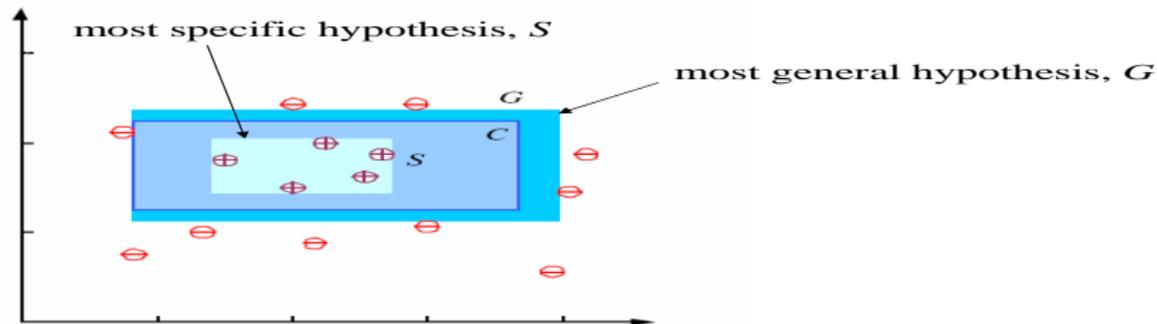


- Generalization error

$$\text{Proba}[\text{error}] = \mathbb{E} [\mathbb{I}(\hat{y}(\mathbf{x}) \neq y)] = \int \mathbb{I}(\hat{y}(\mathbf{x}) \neq y) p(\mathbf{x}, y) d\mathbf{x}dy$$

where the expectation is w.r.t the UNKNOWN probability density function $p(\mathbf{x}, y)$ of (\mathbf{x}, y) .

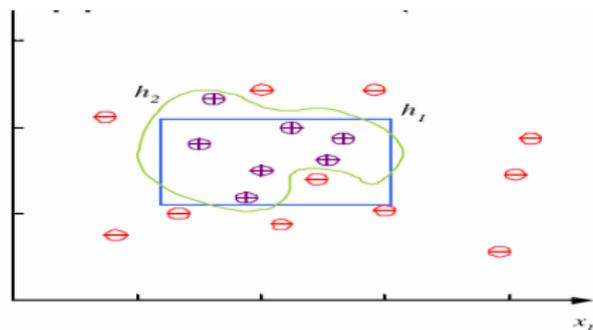
Version Space



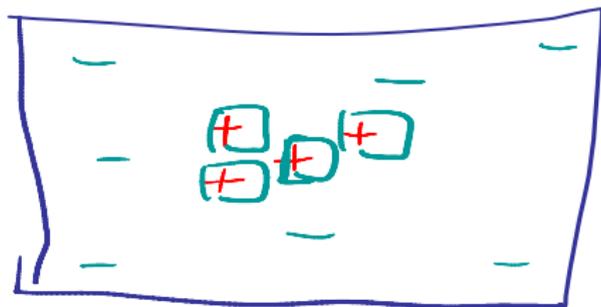
- There may be many functions which have zero training error, ranging from the most specific hypothesis to the most general. Which one we pick depends on our prior knowledge.

Lower Bound on Achievable Error Rate

- If the true concept (green blob) is a rectangle, we can fit it perfectly, and thus get 0 training error. But if the truth is more complex, we will just choose the best-fitting rectangular approximation (blue box) and so $\text{Err}_{\text{Train}} = N_{\text{err}} \neq 0$.



Overfitting in Rectangle Land



- We can always make the empirical/training error be 0 by putting a little rectangle around every + in the training set.
- This may not lead to good generalization performance
- Hence we cannot use directly empirical error to select between models of different complexity.