

Day 6 - Linear Regression and Logistic Regression

Agenda:

- Linear Regression
 - Examples
 - Issues to Pay Attention To with Linear Regression
- Classification and Logistic Regression
 - Training classifiers
 - Evaluating classifiers

More thoughts on square capital example and whether to approach problem as regression or classification

Least Squares Formulation for Linear Regression (for a general model)

Given: $D = \{(X^{(i)}, y_i)\}_{i=1, \dots, n}$, $X^{(i)} \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

Model: $y = \underbrace{\theta_1 g_1(X^{(i)}) + \dots + \theta_k g_k(X^{(i)})}_{f_{\theta}(x)} + \text{Error}$

Want $y \approx \bar{X} \theta$

w/ $\bar{X} = \begin{pmatrix} g_1(X^{(1)}) & g_2(X^{(1)}) & \dots & g_k(X^{(1)}) \\ \vdots & \ddots & \ddots & \vdots \\ g_n(X^{(n)}) & \dots & \dots & g_k(X^{(n)}) \end{pmatrix} = \begin{pmatrix} - & g(X^{(1)}) & - \\ - & g(X^{(2)}) & - \\ \vdots & \vdots & \vdots \\ - & g(X^{(n)}) & - \end{pmatrix}$

Find $\min_{\theta} \frac{1}{2} \|y - \bar{X} \theta\|^2$

where $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$, $\theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_k \end{pmatrix}$

Solution given by Normal Equations

$$X^t X \theta = X^t y \Rightarrow \theta = (X^t X)^{-1} X^t y.$$

Examples of setting up and solving linear regression

Find best fit cubic through 1d data

Data = $\{ (x_i, y_i) \}_{i=1 \dots n}$ w/ $x_i, y_i \in \mathbb{R}$

Model $y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \text{noise}$

Find

$$\min_{\theta} \frac{1}{2} \|y - X\theta\|^2$$

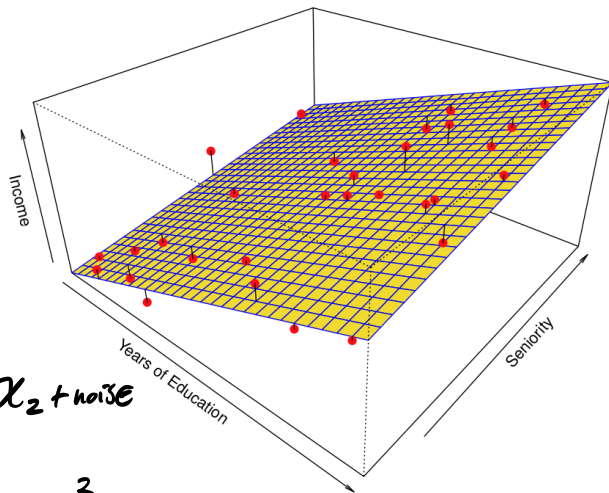
where $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$, $\theta = \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_3 \end{pmatrix}$, $X = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{pmatrix}$

Solution given by Normal Equations

$$X^t X \theta = X^t y \Rightarrow \theta = (X^t X)^{-1} X^t y.$$

Find best fit plane

$$\text{Data} = \left\{ \underbrace{(x^{(i)})}_{\mathbb{R}^2}, \underbrace{y_i}_{\mathbb{R}} \right\}_{i=1 \dots n}$$



Model: $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \text{noise}$

Find

$$\min_{\theta} \frac{1}{2} \|y - \bar{X}\theta\|^2$$

where $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$, $\theta = \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_3 \end{pmatrix}$, $\bar{X} = \begin{pmatrix} \end{pmatrix}$

Solving and Optimization Problem using Gradient Descent

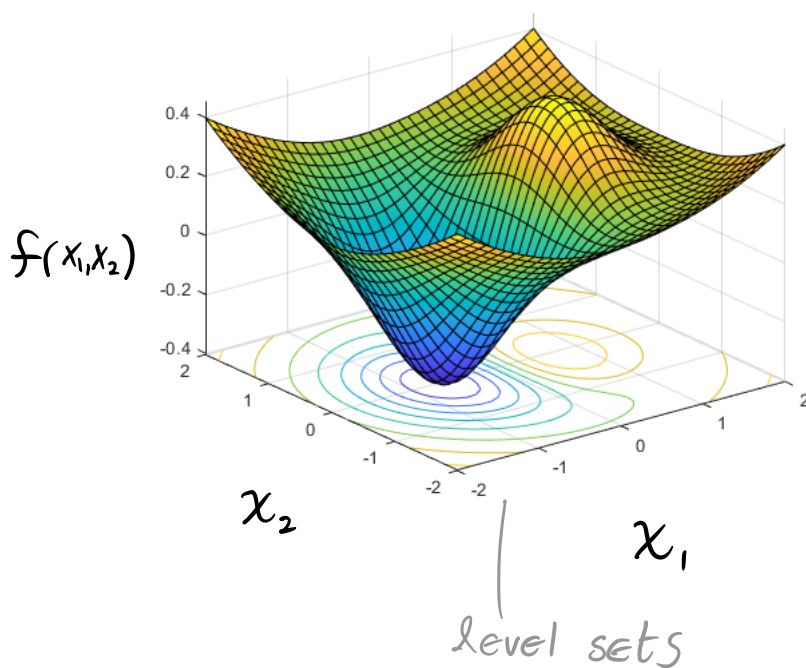
$$\min_{\mathbf{x}} f(\mathbf{x})$$

Gradient descent: Take successive steps "downhill"

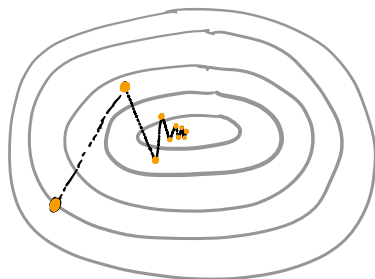
$$\mathbf{x}^{i+1} = \mathbf{x}^i - \alpha \nabla f(\mathbf{x}^i)$$

step size,
learning rate

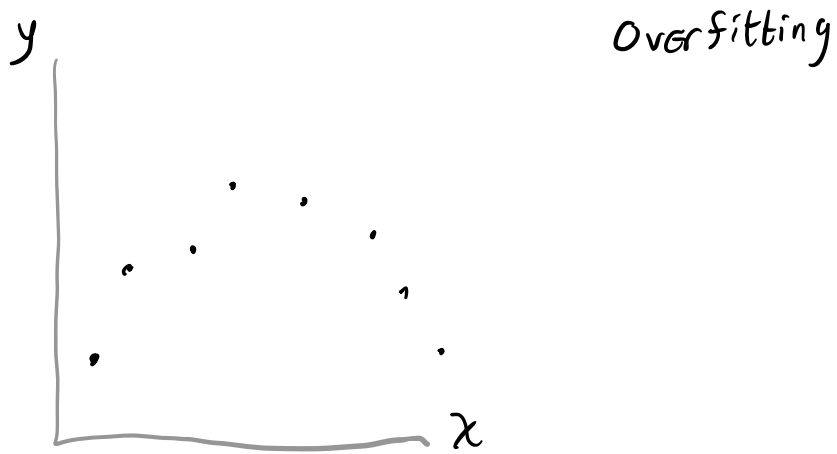
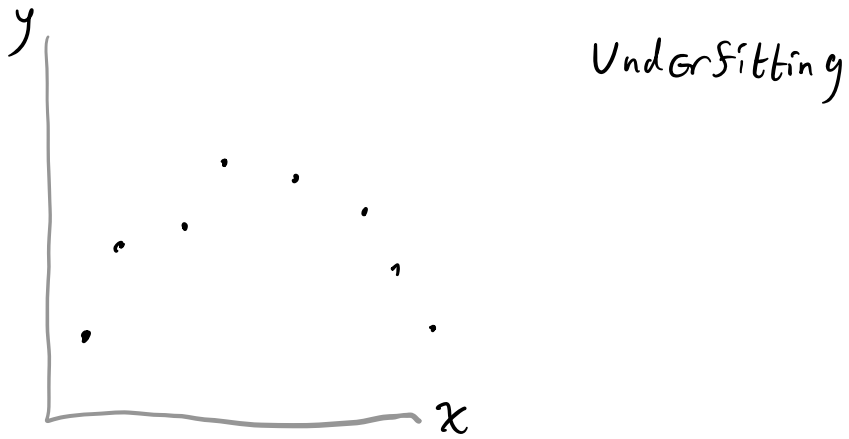
$-\nabla f$ points in direction
of steepest descent



Depiction of gradient descent



Things that can go wrong: Underfitting and Overfitting



Things that can go wrong: numerical instability

Other topics:

What happens when there is fewer data than features?

What happens if there are outliers in the data?

How do you deal with categorical features?

Be careful about whether you want to view your problem as a prediction task

Classification and Logistic Regression

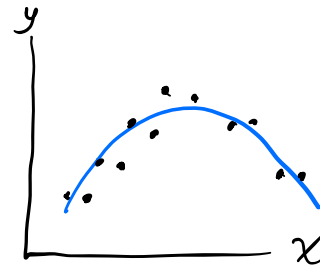
Viewing Regression and Classification as function estimation problems

Regression : predict a continuous value

$$\text{Let } f: \mathbb{R}^d \rightarrow \mathbb{R}$$
$$y = f(x) + \text{noise}$$

$$\text{Given } \circ \{ (x^{(i)}, y_i) \}_{i=1 \dots n}$$

$$\text{Find } \circ f$$



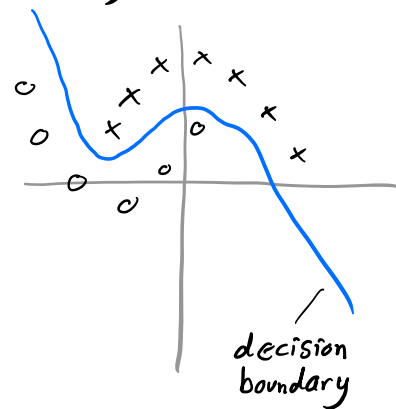
Classification : predict membership in a category

$$\text{Let } f: \mathbb{R}^d \rightarrow \begin{Bmatrix} \text{cat } 1 \\ \vdots \\ \text{cat } m \end{Bmatrix}$$

$$y = f(x) + \text{noise}$$

$$\text{Given } \circ \{ (x^{(i)}, y_i) \}_{i=1 \dots n}$$

$$\text{Find } \circ f$$



Terminology :

- x - input variables, predictors, independent vars, features
- y - response, dependent variable, output variable
- f - model, predictor, hypothesis

Parametric Approach: Choose a model for f with unknown parameters. Estimate the parameters.

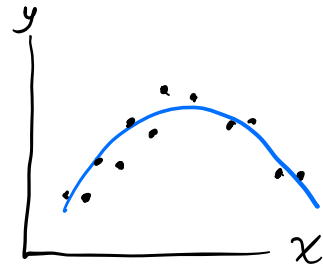
Parametric

Regression: predict a continuous value

$$\text{Model } f_{\theta}: \mathbb{R}^d \rightarrow \mathbb{R}$$
$$y = f_{\theta}(x) + \text{noise}$$

$$\text{Given: } \{(x^{(i)}, y_i)\}_{i=1, \dots, n}$$

Find: θ



Parametric

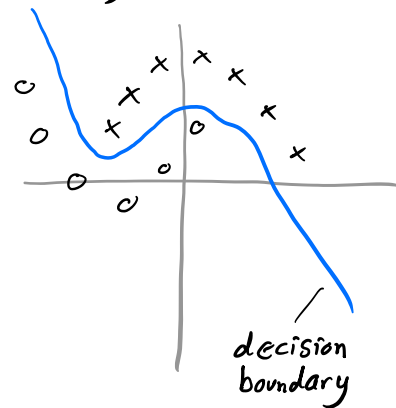
Classification: predict membership in a category

$$\text{Let } f_{\theta}: \mathbb{R}^d \rightarrow \begin{cases} \text{cat 1} \\ \vdots \\ \text{cat m} \end{cases}$$

$$y = f_{\theta}(x) + \text{noise}$$

$$\text{Given: } \{(x^{(i)}, y_i)\}_{i=1, \dots, n}$$

Find: θ



Approach for estimating θ :

Select a model for f w/ parameters θ

&

minimize the loss between
training labels and predictions on
training data

$$\min_{\theta} \sum_{i=1}^n L(y_i, f_{\theta}(x^{(i)}))$$

|
loss
function

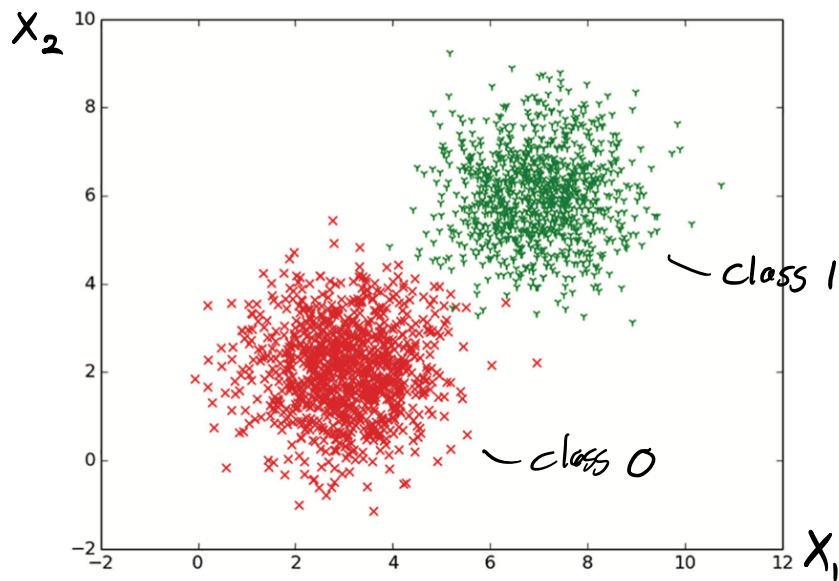
Example: linear regression $L(y, \hat{y}) = |y - \hat{y}|^2$
 "square loss" or ℓ^2 loss

Binary Classification in 2D with logistic regression

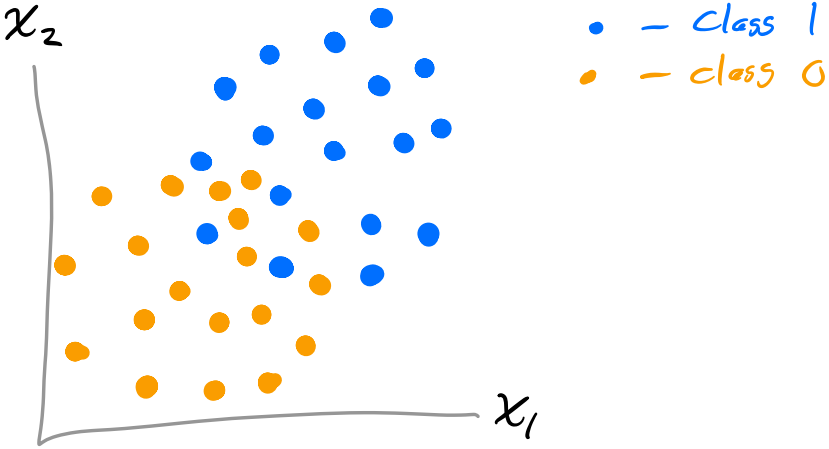
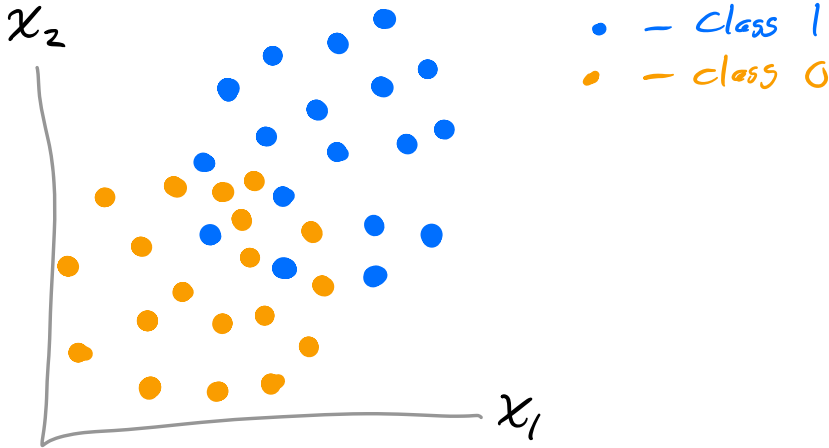
Training Data $\{x^{(i)}, y_i\}_{i=1 \dots n}$

\mathbb{R}^2 \mathbb{R}

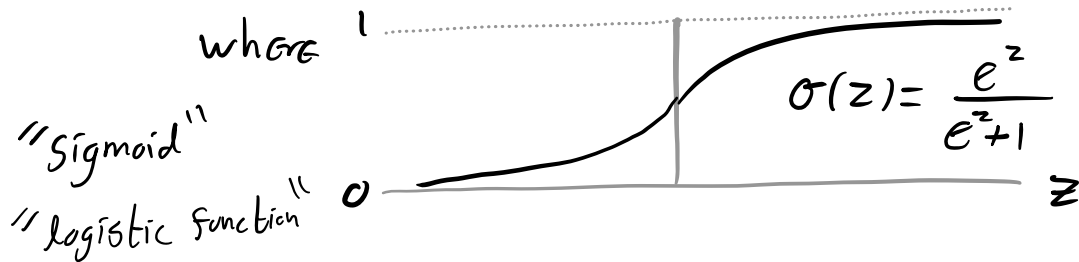
$y_i = \begin{cases} 1 & \text{if class 1} \\ 0 & \text{if class 0} \end{cases}$



Given this data, draw a decision boundary (curve where you would say class 1 is on one side and class 2 is on the other side)



Model $y = \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = \hat{y}(x; \theta)$



Solve $\min_{\theta} \sum_{i=1}^n L(y_i, \hat{y}(x^{(i)}; \theta))$ for $\hat{\theta}$

Predict: For new sample x , predict

$$\begin{cases} \text{class 1} & \text{if } \hat{y} \geq \frac{1}{2} \\ \text{class 0} & \text{if } \hat{y} < \frac{1}{2} \end{cases}$$

What loss function should you use?

One choice - log loss

$$\begin{aligned} L(y, \hat{y}) &= \begin{cases} -\log(\hat{y}) & \text{if } y=1 \\ -\log(1-\hat{y}) & \text{if } y=0 \end{cases} \\ &= -y \log \hat{y} - (1-y) \log(1-\hat{y}) \end{aligned}$$

Decision Boundary for Logistic Regression

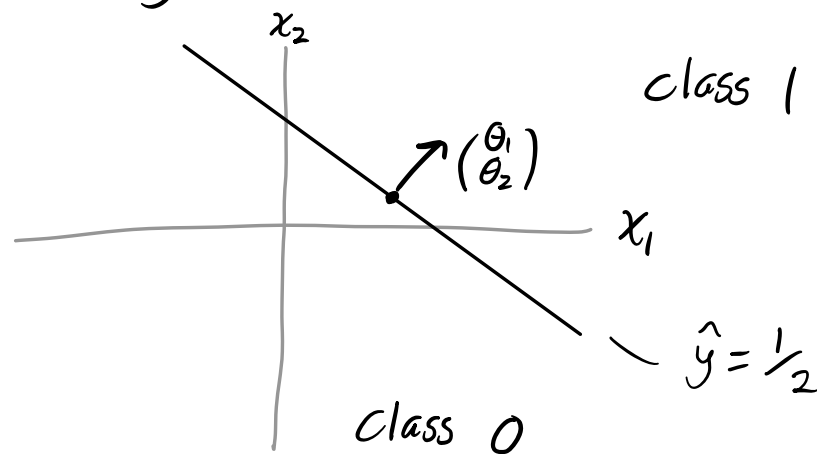
Training Data: $\{x^{(i)}, y_i\}_{i=1, \dots, n}$ $\begin{matrix} \mathbb{R}^2 \\ \mathbb{R} \end{matrix}$ $y_i = \begin{cases} 1 & \text{if class 1} \\ 0 & \text{if class 0} \end{cases}$

Model: $y = \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = \hat{y}(x; \theta)$

Predict: For new sample x , predict

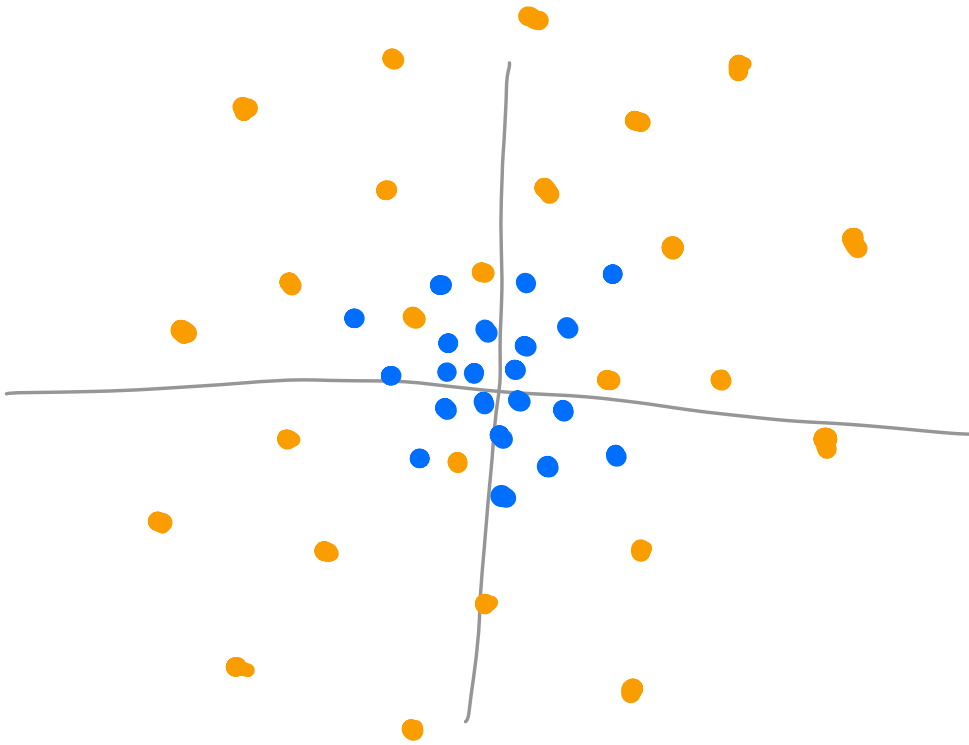
$$\begin{cases} \text{class 1} & \text{if } \hat{y} \geq \frac{1}{2} \\ \text{class 0} & \text{if } \hat{y} < \frac{1}{2} \end{cases}$$

Decision boundary is linear



Activity:

Could you use logistic regression to build a reasonable classifier for the following data?



Evaluating Classifiers

Prediction

		+	-
Truth	+	True Positive (TP)	False Negative (FN)
	-	False Positive (FP)	True Negative (TN)

$$\text{Accuracy} \% = \frac{TP + TN}{\text{total}}$$

Activity: Someone invents a test for a rare disease that affects 0.1% of the population. The test has accuracy 99.9%. Are you convinced this is a good test?

$$\text{Precision} \% = \frac{TP}{TP+FP}$$

What fraction of predicted positives are real?

$$\text{Recall} \% = \frac{TP}{TP+FN}$$

What fraction of positives are correctly predicted?

Want high precision & high recall

Activity: You are building a binary classifier that detects whether a pedestrian is crossing the sidewalk within 30 feet of a self driving car. If the detection is positive, the car puts on the breaks. Would you rather have good precision and great recall or good recall and great precision?

There is a trade off between True Positives and False Positives, and between True Negatives and False Negatives

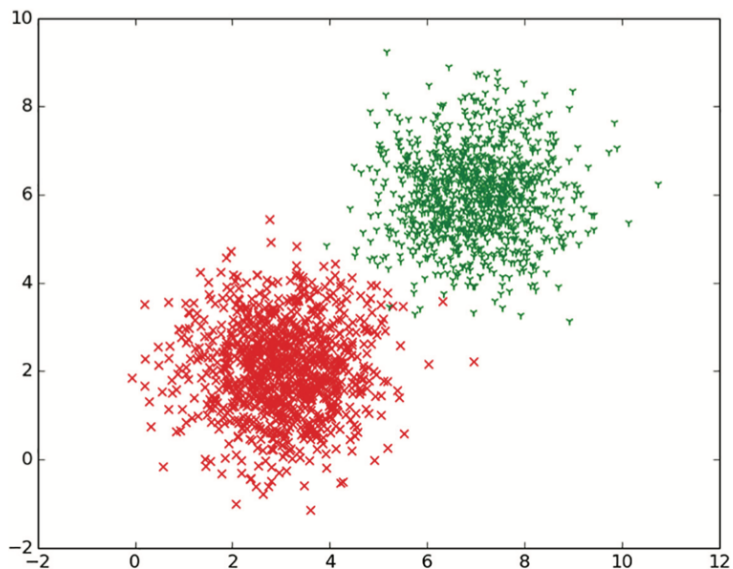
$$\text{Training Data: } \left\{ \overset{\mathbb{R}^2}{x^{(i)}}, \overset{\mathbb{R}}{y_i} \right\}_{i=1 \dots n} \quad y_i = \begin{cases} 1 & \text{if class 1} \\ 0 & \text{if class 0} \end{cases}$$

$$\text{Model: } y = \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = \hat{y}(x; \theta)$$

Predict: For new sample x , predict

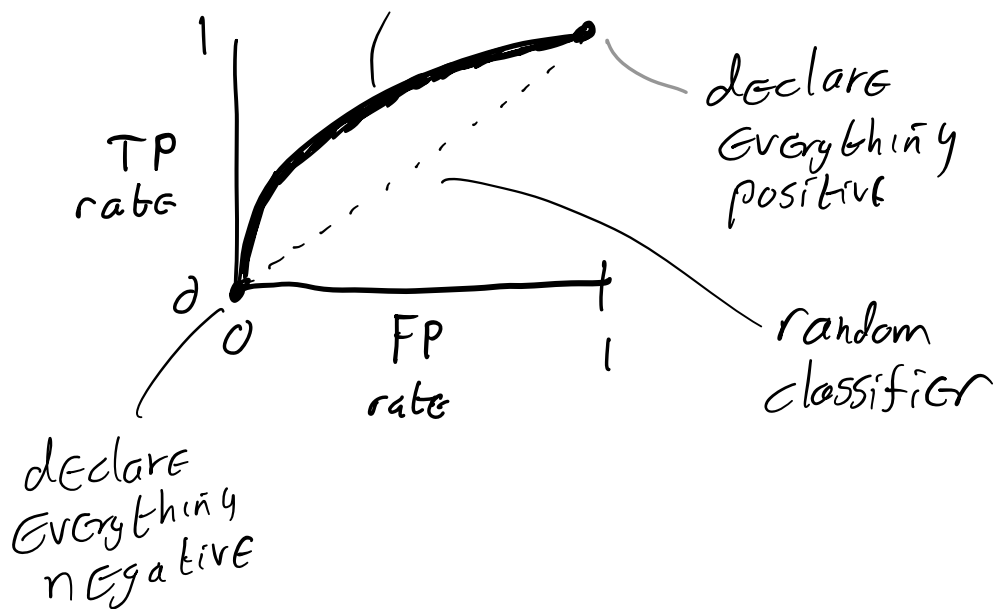
$$\begin{cases} \text{class 1} & \text{if } \hat{y} \geq \frac{1}{2} \\ \text{class 0} & \text{if } \hat{y} < \frac{1}{2} \end{cases}$$

~ could choose any value



Receiver Operating Characteristic Curves

Each point is a classifier w/ different threshold



Comparing classifiers and Area-Under-Curve (AUC)

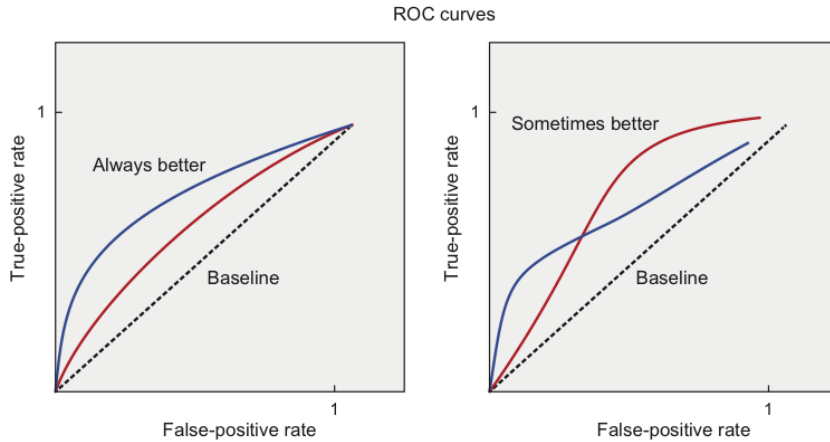
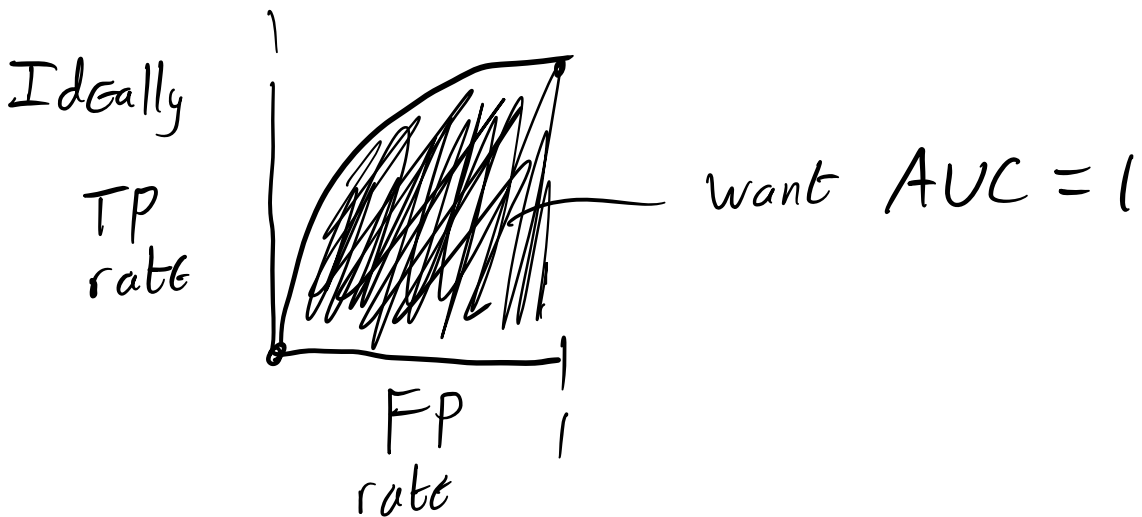


Figure 5.6 The principled way to compare algorithms is to examine their ROC curves. When the true-positive rate is greater than the false-positive rate in every situation, it's straightforward to declare that one algorithm is dominant in terms of its performance. If the true-positive rate is less than the false-positive rate, the plot dips below the baseline shown by the dotted line.



Also common to plot precision-recall curves

