

10-701: Introduction to Machine Learning Lecture 7 – Logistic Regression

Pradeep Ravikumar

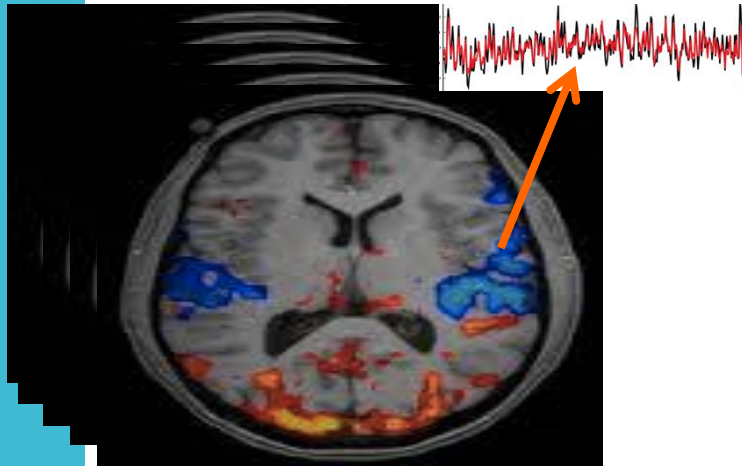
Spring 2026

Front Matter

- Announcements:
 - Recitation this Friday
- Recommended Readings:
 - Murphy, [Section 8.1 - 8.3](#)

Regression to Classification

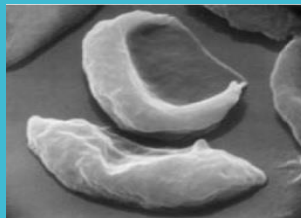
Regression



Y = Age of a subject

X = Brain Scan

Classification



Anemic cell
Healthy cell

X = Cell Image

Y = Diagnosis

Can we predict the “probability” of class label being Anemic or Healthy – a real number – using regression methods?

But output (probability) needs to be in $[0,1]$

Classification

Goal: Construct a **predictor** $f : X \rightarrow Y$ to minimize a risk (performance measure) $R(f)$



Features, X



Sports
Science
News

Labels, Y

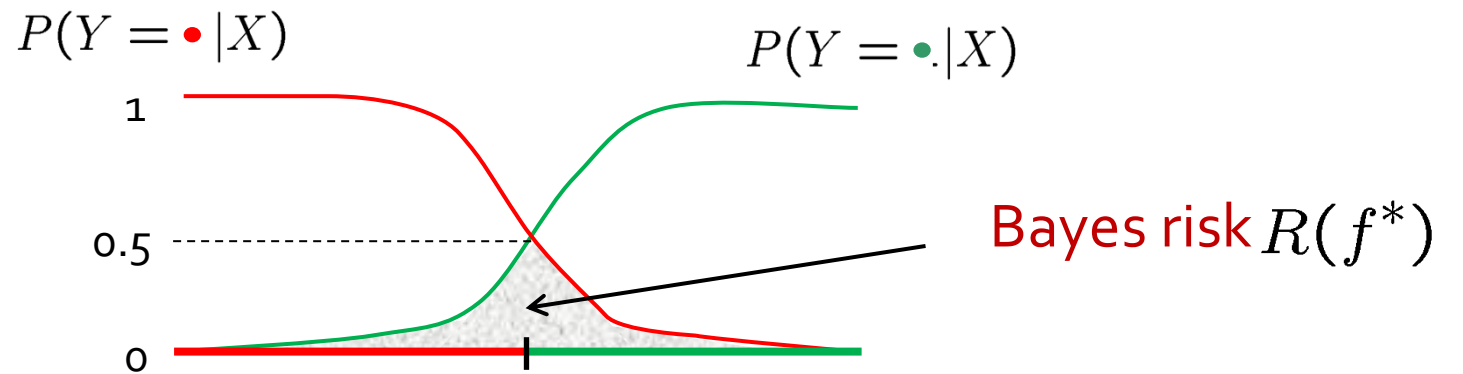
$$R(f) = P(f(X) \neq Y)$$

Probability of Error

Optimal Classification

Optimal predictor:
(Bayes classifier)

$$f^* = \arg \min_f P(f(X) \neq Y)$$



$$f^*(x) = \arg \max_{Y=y} P(Y = y | X = x)$$

- Even the optimal classifier makes mistakes $R(f^*) > 0$
- Optimal classifier depends on **unknown** distribution P_{XY}

Optimal Classifier

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Bayes Rule:

$$P(Y = y|X = x) = \frac{P(X = x|Y = y)P(Y = y)}{P(X = x)}$$

$$f^*(x) = \arg \max_{Y=y} P(Y = y|X = x)$$

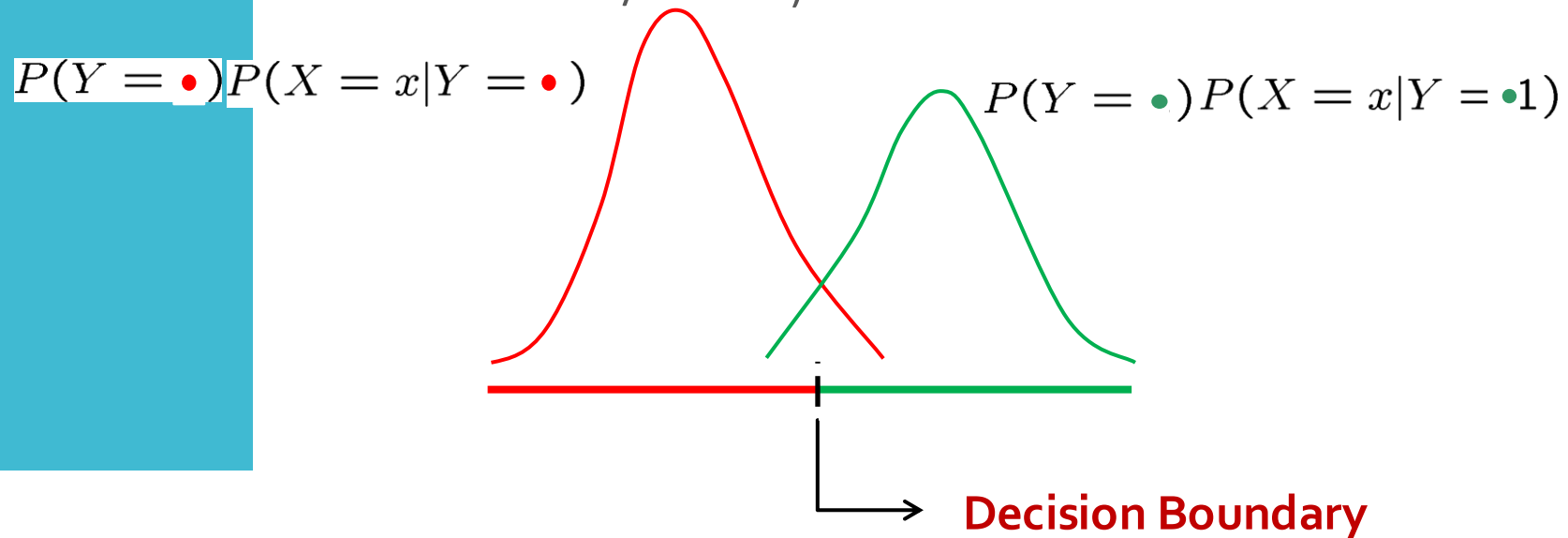
$$= \arg \max_{Y=y} \underbrace{P(X = x|Y = y)}_{\text{Class conditional density}} \underbrace{P(Y = y)}_{\text{Class prior}}$$

Class conditional density Class prior

Example Decision Boundaries

$$P(X = x|Y = y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x - \mu_y)^2}{2\sigma_y^2}\right)$$

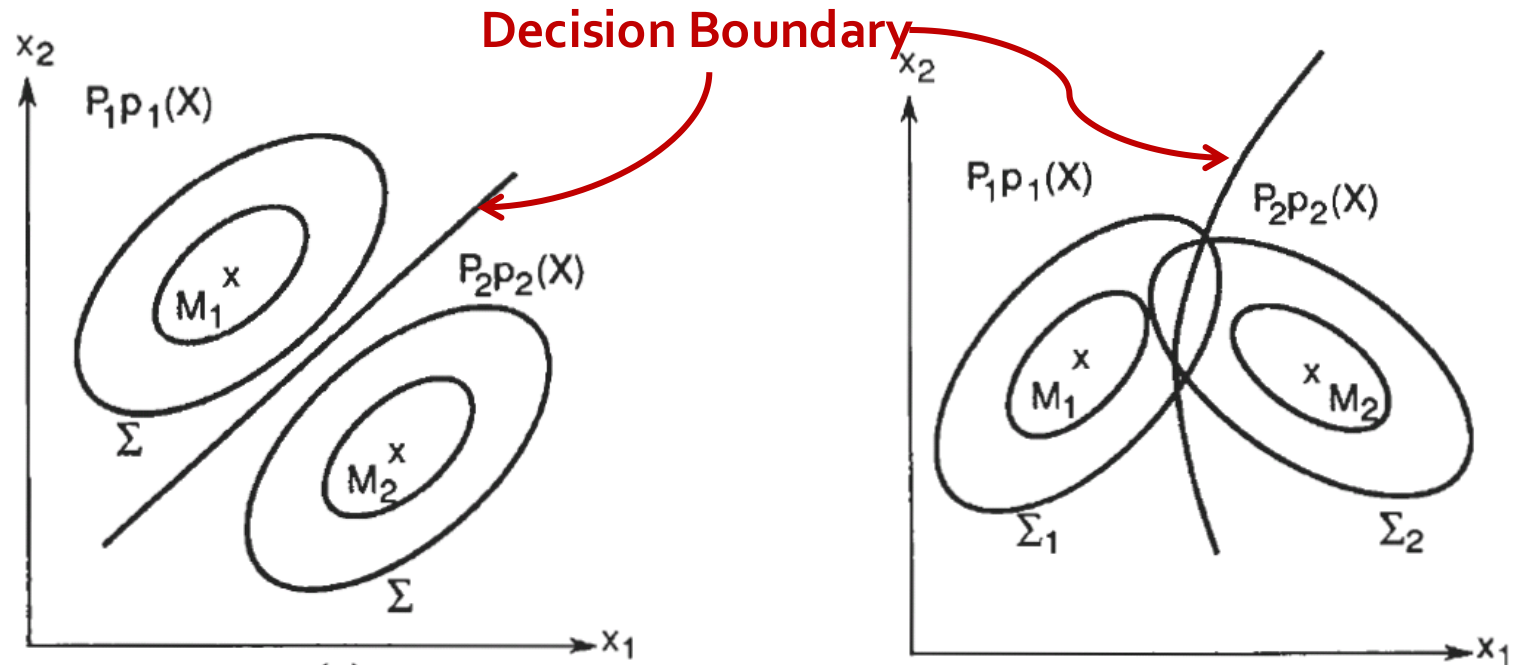
- Gaussian class conditional densities (1-dimension/feature)



Example Decision Boundaries

- Gaussian class conditional densities (2-dimensions/features)

$$P(X = x|Y = y) = \frac{1}{\sqrt{2\pi|\Sigma_y|}} \exp\left(-\frac{(x - \mu_y)\Sigma_y^{-1}(x - \mu_y)'}{2}\right)$$



Learning the Optimal Classifier

Optimal classifier:

$$\begin{aligned} f^*(x) &= \arg \max_{Y=y} P(Y = y | X = x) \\ &= \arg \max_{Y=y} \underbrace{P(X = x | Y = y)}_{\text{Class conditional density}} \underbrace{P(Y = y)}_{\text{Class prior}} \end{aligned}$$

Class conditional
density Class prior

Need to know Prior $P(Y = y)$ for all y

Likelihood $P(X=x|Y = y)$ for all x, y

Two truths and a lie

Which statement is the LIE?

- A. Bayes classifier: $f^*(x) = \operatorname{argmax}_y P(Y = y | X = x)$
- B. Even the optimal classifier can make mistakes (Bayes risk $R(f^*) > 0$)
- C. The Bayes classifier does NOT depend on the (unknown) distribution of (X, Y)

Discuss with a neighbor for 20 seconds.

Generative vs. Discriminative Classifiers

Generative classifiers

- Assume some functional form for $P(X,Y)$ (or $P(X|Y)$ and $P(Y)$)
- Estimate parameters of $P(X|Y)$, $P(Y)$ directly from training data
- Use Bayes rule to calculate $P(Y|X)$

Discriminative classifiers

- Assume some functional form for $P(Y|X)$ or for the decision boundary
- Estimate parameters of $P(Y|X)$ directly from training data

Probabilistic Classification

- Classification:
 - (Unknown) Target function, $c^*: \mathcal{X} \rightarrow \mathcal{Y}$
 - Classifier, $h : \mathcal{X} \rightarrow \mathcal{Y}$
 - Goal: find a classifier, h , that best approximates c^*
- Probabilistic Classification:
 - (Unknown) Target *distribution*, $y \sim P^*(Y|\mathbf{x})$
 - Distribution, $P(Y|\mathbf{x})$
 - Goal: find a distribution, P , that best approximates P^*

Building a Probabilistic Classifier

- Define a decision rule
 - Given a test data point \mathbf{x}' , predict its label \hat{y} using the *posterior distribution* $P(Y = y|X = \mathbf{x}')$
 - Common choice: $\hat{y} = \underset{y}{\operatorname{argmax}} P(Y = y|X = \mathbf{x}')$
- Model the posterior distribution
 - Option 1 (“Discriminative Model”) (today)
 - Model $P(Y|X)$ directly as some function of X
 - Option 2 (“Generative Model”)
 - Use Bayes’ rule:

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \propto P(X|Y) P(Y)$$

Modelling the Posterior

- Suppose we have binary labels $y \in \{0,1\}$ and D -dimensional inputs $\mathbf{x} = [1, x_1, \dots, x_D]^T \in \mathbb{R}^{D+1}$

- **Assume**

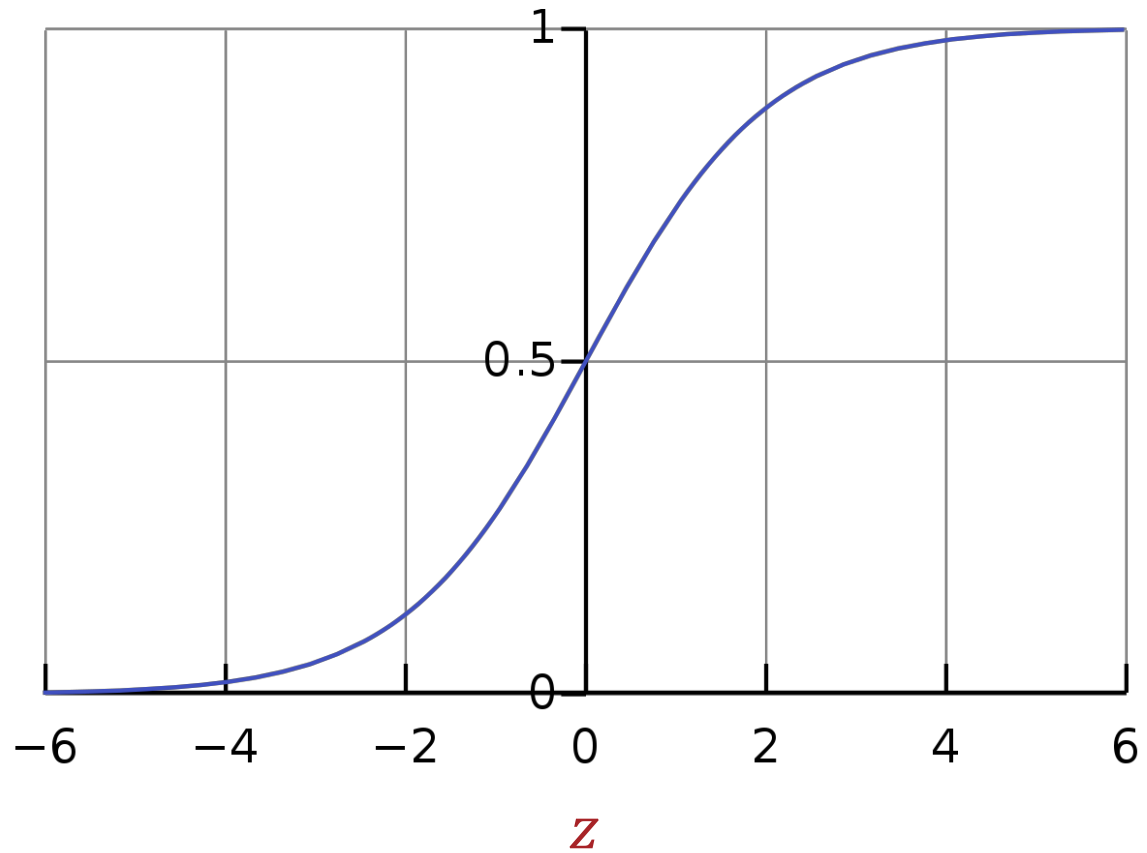
$$\begin{aligned} P(Y = 1|\mathbf{x}) = \text{logit}(\mathbf{w}^T \mathbf{x}) &= \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \\ &= \frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1} \end{aligned}$$

- This implies two useful facts:

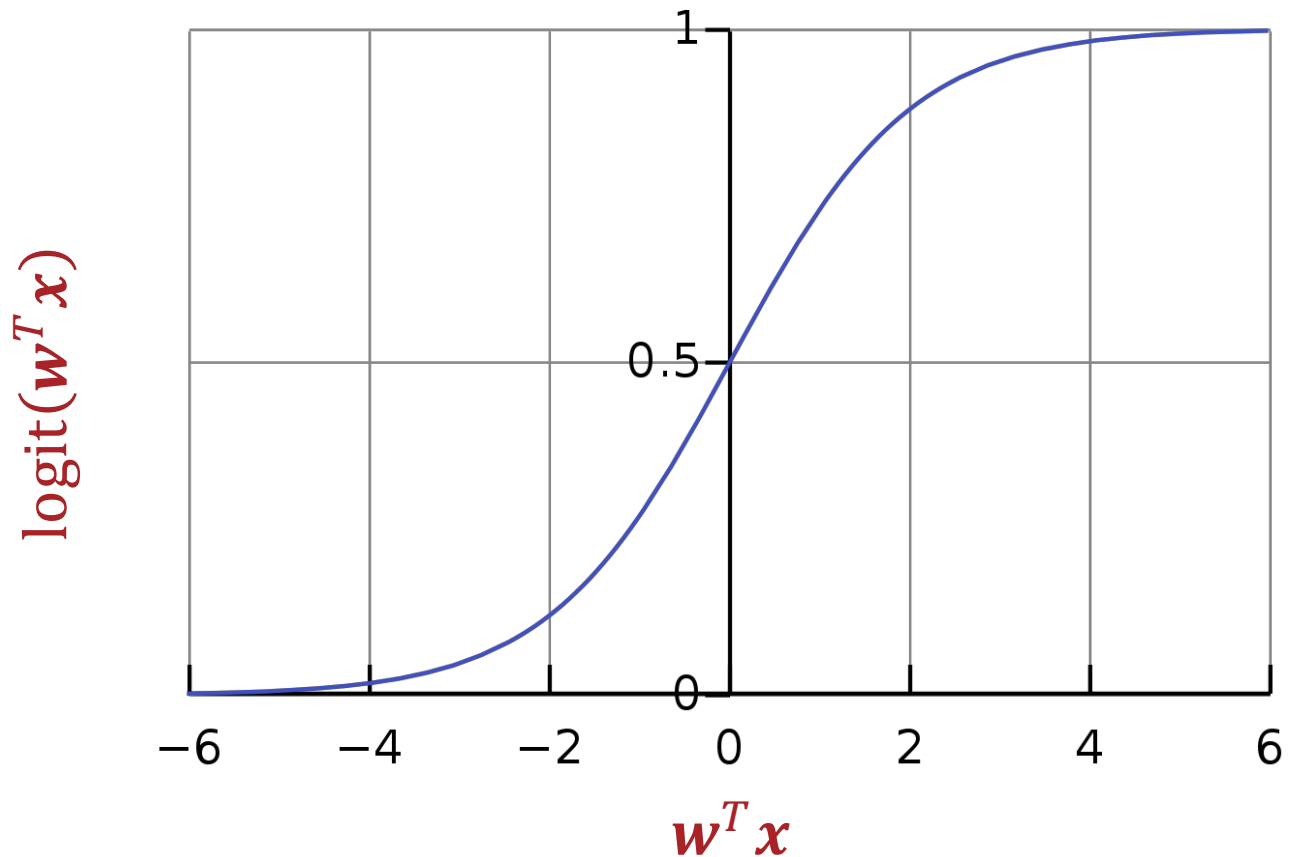
1. $P(Y = 0|\mathbf{x}) = 1 - P(Y = 1|\mathbf{x}) = \frac{1}{\exp(\mathbf{w}^T \mathbf{x}) + 1}$
2. $\frac{P(Y = 1|\mathbf{x})}{P(Y = 0|\mathbf{x})} = \exp(\mathbf{w}^T \mathbf{x}) \rightarrow \log \frac{P(Y = 1|\mathbf{x})}{P(Y = 0|\mathbf{x})} = \mathbf{w}^T \mathbf{x}$

Logistic Function

$$\text{logit}(z) = \frac{1}{1 + e^{-z}}$$



Why use the Logistic Function?



- Differentiable everywhere
- $\text{logit}: \mathbb{R} \rightarrow [0, 1]$
- The decision boundary is linear in \mathbf{x} !

Mythbuster

Mythbusters

MYTH

Because the sigmoid is nonlinear, logistic regression has a nonlinear decision boundary.

Mythbuster

Mythbusters

MYTH

Because the sigmoid is nonlinear, logistic regression has a nonlinear decision boundary.

BUSTED

With threshold $1/2$, the boundary is $w^T x = 0$ — linear in x . (The sigmoid only squashes scores into probabilities.)

Logistic Regression Decision Boundary

$$\hat{y} = \begin{cases} 1 & \text{if } P(Y = 1|\mathbf{x}) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

$$P(Y = 1|\mathbf{x}) = \text{logit}(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \geq \frac{1}{2}$$

$$2 \geq 1 + \exp(-\mathbf{w}^T \mathbf{x})$$

$$1 \geq \exp(-\mathbf{w}^T \mathbf{x})$$

$$\log(1) \geq -\mathbf{w}^T \mathbf{x}$$

$$0 \leq \mathbf{w}^T \mathbf{x}$$

Logistic Regression Decision Boundary

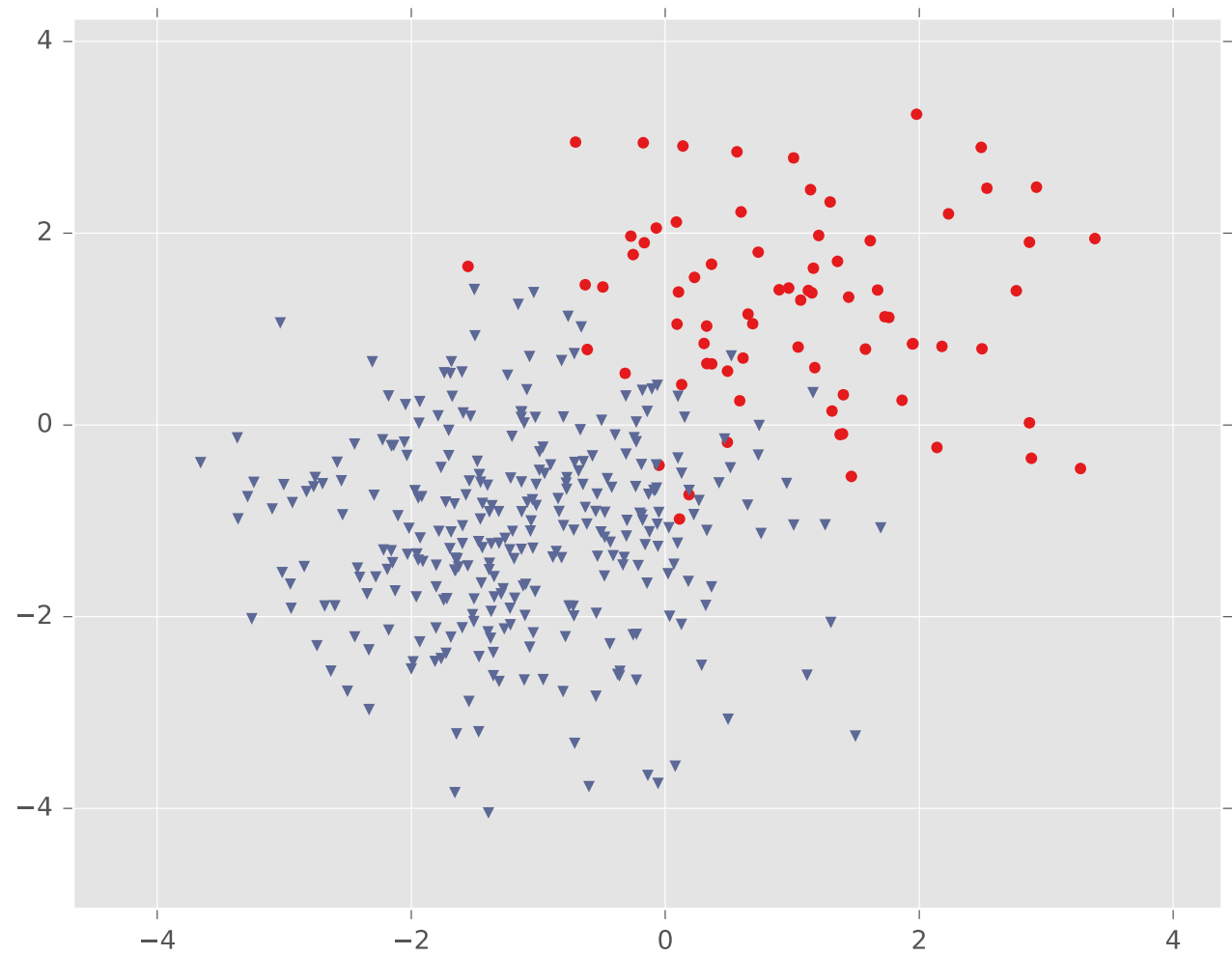


Figure courtesy of Matt Gormley

Logistic Regression Decision Boundary

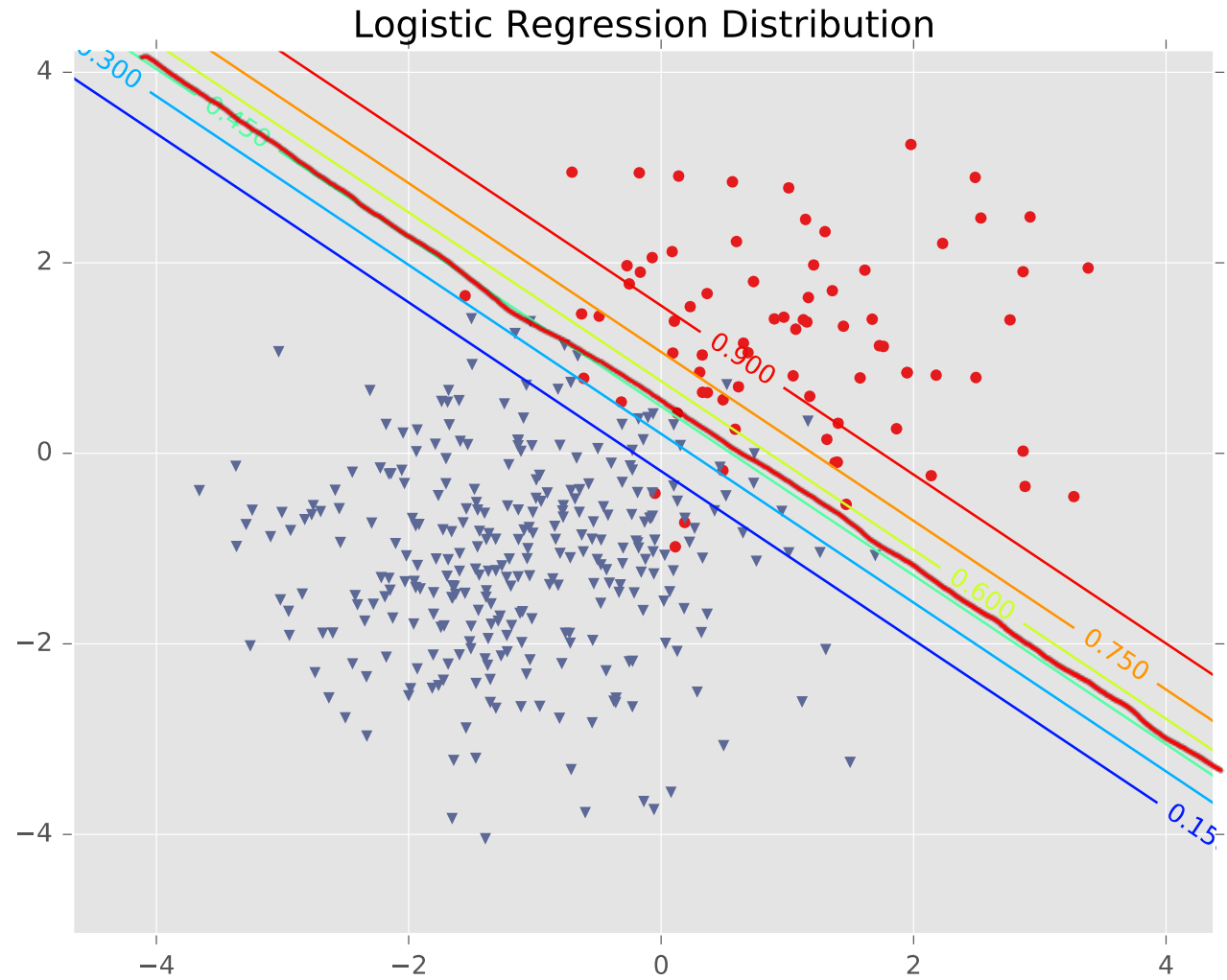


Figure courtesy of Matt Gormley

Logistic Regression Decision Boundary



Figure courtesy of Matt Gormley

Mythbuster

Mythbusters

MYTH

“Logistic regression” predicts a real-valued target like linear regression.

Mythbuster

Mythbusters

MYTH

“Logistic regression” predicts a real-valued target like linear regression.

BUSTED

It models $P(Y = 1 | X)$ using a logistic (sigmoid) of $w^T x$, producing a probability in $[0,1]$; we classify by thresholding.

Logistic Regression

Not really regression

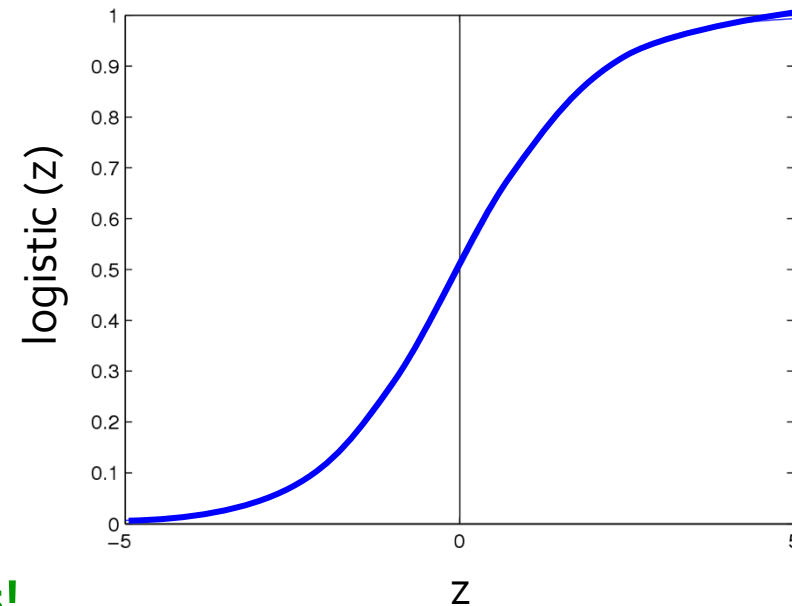
Assumes the following functional form for $P(Y|X)$:

$$P(Y = 0|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Logistic function applied to a linear function of the data

Logistic function
(or Sigmoid):

$$\frac{1}{1 + \exp(-z)}$$



Features can be discrete or continuous!

Logistic Regression is a Linear Classifier!

Assumes the following functional form for $P(Y|X)$:

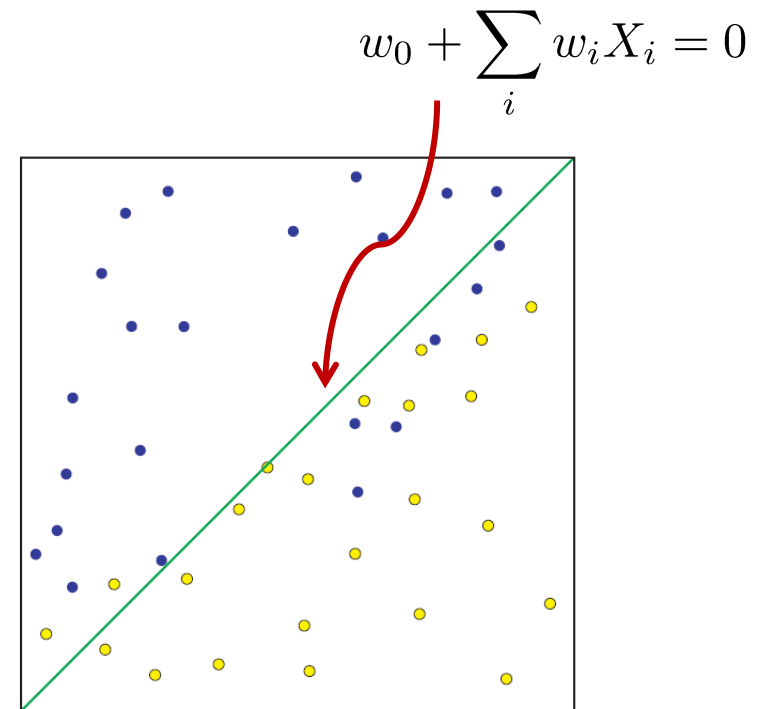
$$P(Y = 0|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Decision boundary: Note - Labels are 0,1

$$P(Y = 0|X) \underset{1}{\geq} P(Y = 1|X)$$

$$w_0 + \sum_i w_i X_i \underset{0}{\geq} 1$$

(Linear Decision Boundary)



General Recipe for Machine Learning

- Define a model and model parameters
- Write down an objective function
- Optimize the objective w.r.t. the model parameters

Recipe for Logistic Regression

- Define a model and model parameters
 - Assume independent, identically distributed (iid) data
 - Assume $P(Y = 1|X) = \text{logit}(\mathbf{w}^T \mathbf{x})$
 - Parameters: $\mathbf{w} = [w_0, w_1, \dots, w_D]$
- Write down an objective function
 - ~~Maximize the *conditional* log-likelihood~~
 - Minimize the negative conditional log-likelihood
- Optimize the objective w.r.t. the model parameters
 - ???

Setting the Parameters via Minimum Negative Conditional (log-)Likelihood Estimation (MCLE)

Find \mathbf{w} that minimizes

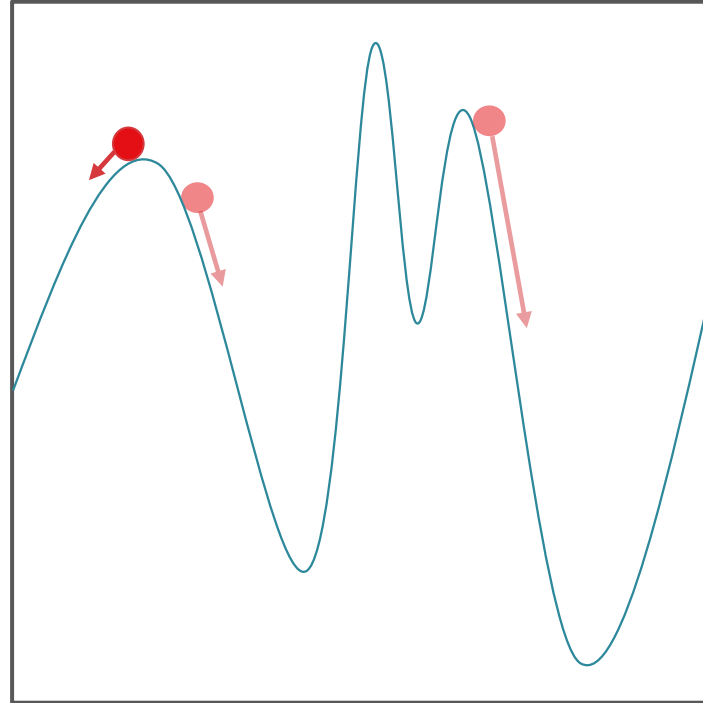
$$\begin{aligned}\ell_{\mathcal{D}}(\mathbf{w}) &= -\log P(y^{(1)}, \dots, y^{(N)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}, \mathbf{w}) = -\log \prod_{n=1}^N P(y^{(n)} | \mathbf{x}^{(n)}, \mathbf{w}) \\ &= -\log \prod_{n=1}^N P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w})^{y^{(n)}} \left(P(Y = 0 | \mathbf{x}^{(n)}, \mathbf{w}) \right)^{1-y^{(n)}} \\ &= -\sum_{i=1}^N y^{(i)} \log P(Y = 1 | \mathbf{x}^{(i)}, \mathbf{w}) + (1 - y^{(i)}) \log P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \\ &= -\sum_{i=1}^N y^{(i)} \log \frac{P(Y = 1 | \mathbf{x}^{(i)}, \mathbf{w})}{P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w})} + \log P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \\ &= -\sum_{i=1}^N y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} - \log \left(1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)}) \right)\end{aligned}$$

Minimizing the Negative Conditional (log-)Likelihood

$$\begin{aligned}\ell_{\mathcal{D}}(\mathbf{w}) &= - \sum_{n=1}^N y^{(n)} \mathbf{w}^T \mathbf{x}^{(n)} - \log \left(1 + \exp(\mathbf{w}^T \mathbf{x}^{(n)}) \right) \\ \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}) &= - \sum_{n=1}^N y^{(n)} \nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{x}^{(n)} - \nabla_{\mathbf{w}} \log \left(1 + \exp(\mathbf{w}^T \mathbf{x}^{(n)}) \right) \\ &= - \sum_{n=1}^N y^{(n)} \mathbf{x}^{(n)} - \frac{\exp(\mathbf{w}^T \mathbf{x}^{(n)})}{1 + \exp(\mathbf{w}^T \mathbf{x}^{(n)})} \mathbf{x}^{(n)} \\ &= \sum_{n=1}^N \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}) - y^{(n)})\end{aligned}$$

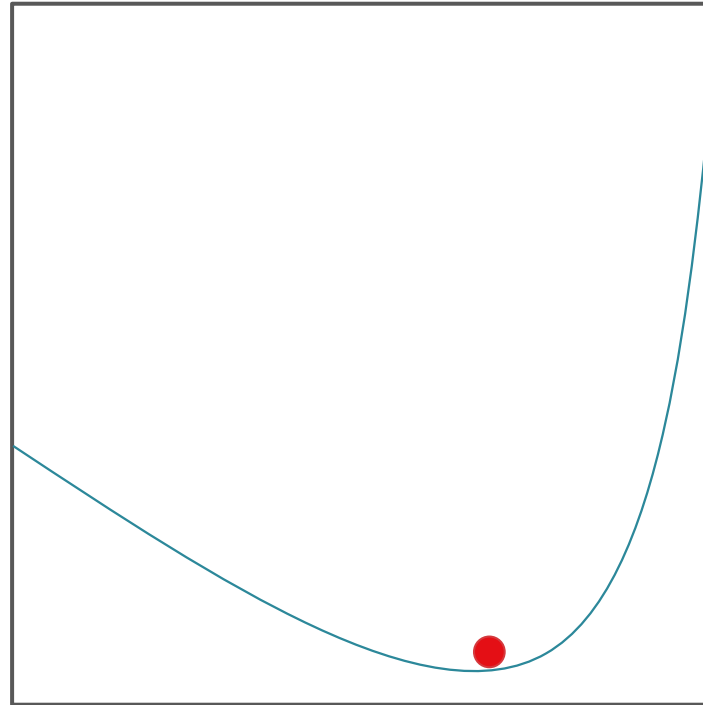
Recall: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



Recall: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



- Good news: the negative conditional log-likelihood, like the squared error, is also convex!

Gradient Descent

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta^{(0)}$

1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$

2. While TERMINATION CRITERION is not satisfied

- a. Compute the gradient:

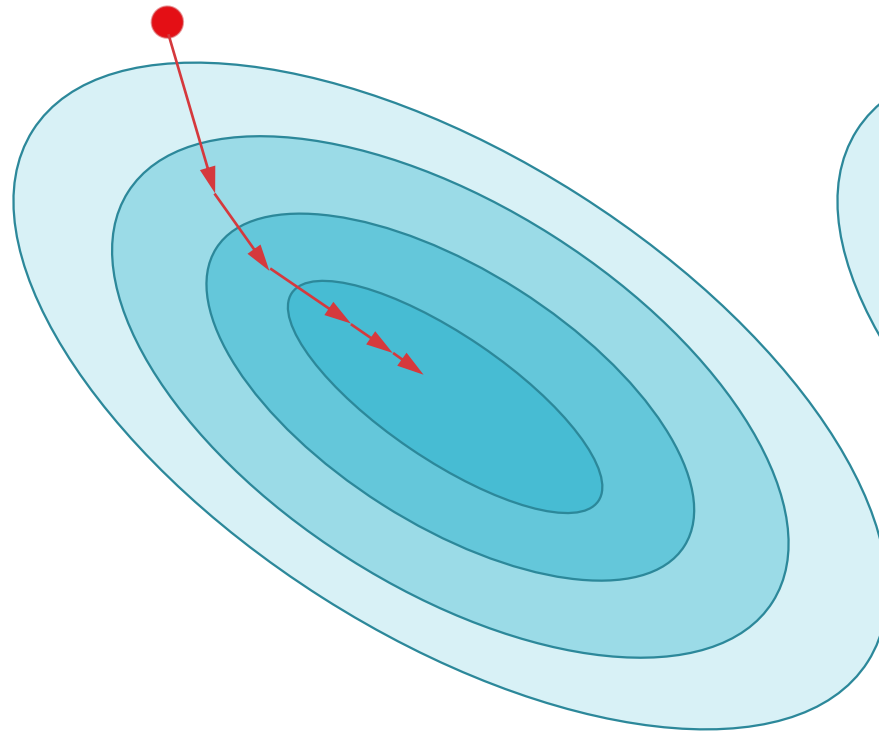
$$O(ND) \left\{ \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)}) = \sum_{n=1}^N \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)}) \right.$$

- b. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$

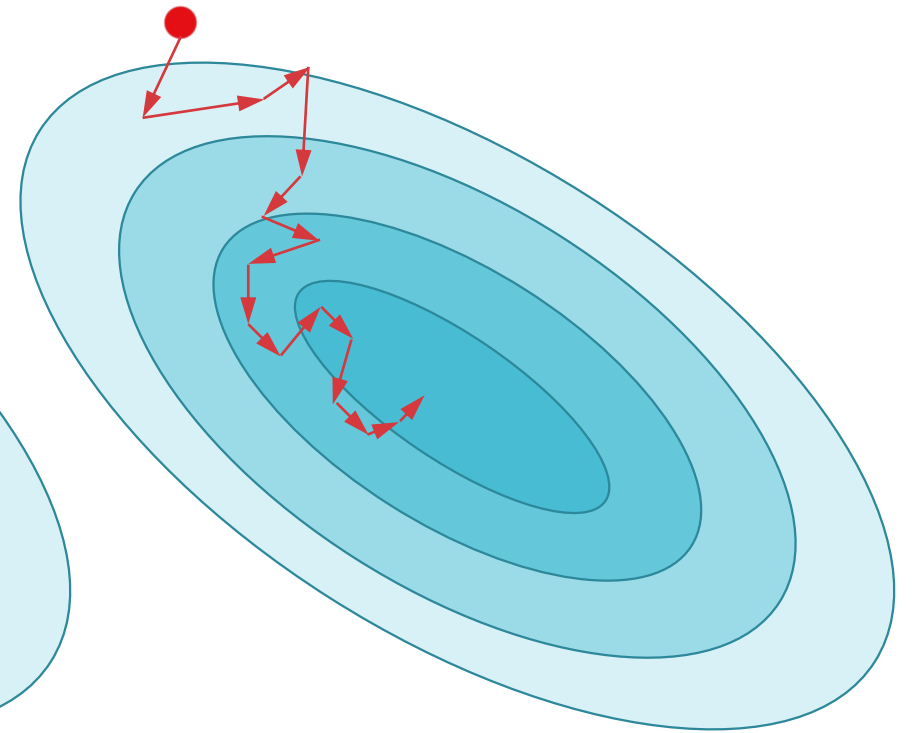
- c. Increment t : $t \leftarrow t + 1$

- Output: $\mathbf{w}^{(t)}$

Stochastic Gradient Descent vs. Gradient Descent



Gradient Descent



Stochastic Gradient Descent

Stochastic Gradient Descent

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{SGD}^{(0)}$
 1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample a data point from \mathcal{D} , $(\mathbf{x}^{(n)}, y^{(n)})$
 - b. Compute the pointwise gradient:
$$\nabla_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)}) = \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)})$$
 - c. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta_{SGD}^{(0)} \nabla_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)})$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{w}^{(t)}$

Stochastic Gradient Descent

- If the data point is sampled uniformly at random, then the expected value of the pointwise gradient is proportional to the full gradient:

$$\begin{aligned} E \left[\nabla_{\mathbf{w}} \ell_{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}}(\mathbf{w}^{(t)}) \right] &= \frac{1}{N} \sum_{n=1}^N \nabla_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)}) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)}) \\ &= \frac{1}{N} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)}) \end{aligned}$$

- In practice, the data set is randomly shuffled then looped through so that each data point is used equally often

Mini-batch Stochastic Gradient Descent

• Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{MB}^{(0)}, B$

1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
2. While TERMINATION CRITERION is not satisfied

a. Randomly sample B data points from \mathcal{D} :

$$\mathcal{D}_{batch} \{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$$

b. Compute the gradient w.r.t. the sampled *batch*:

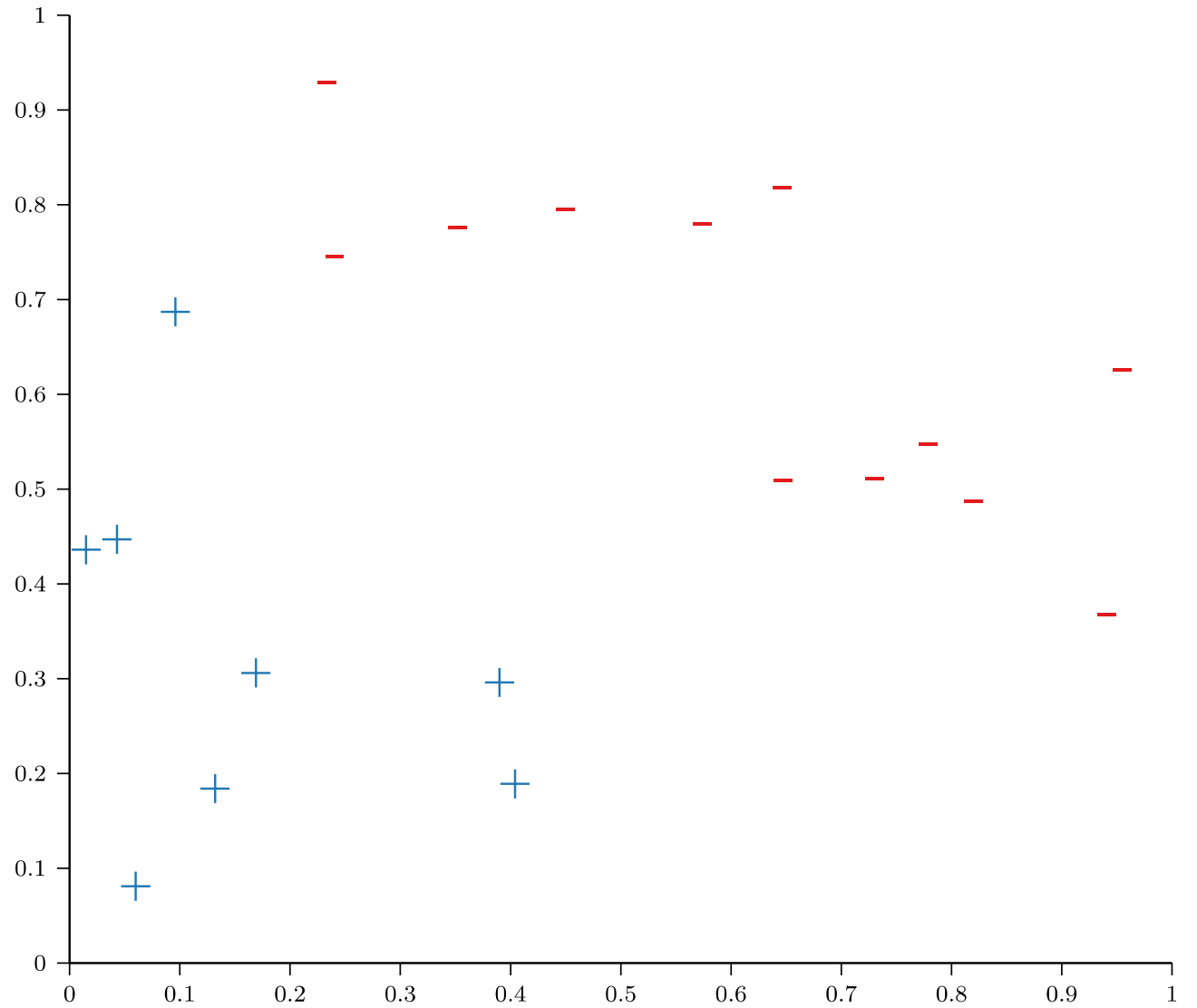
$$\nabla_{\mathbf{w}} \ell_{\mathcal{D}_{batch}}(\mathbf{w}^{(t)}) = \sum_{b=1}^B \mathbf{x}^{(b)} (P(Y = 1 | \mathbf{x}^{(b)}, \mathbf{w}) - y^{(b)})$$

c. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta_{MB}^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}_{batch}}(\mathbf{w}^{(t)})$

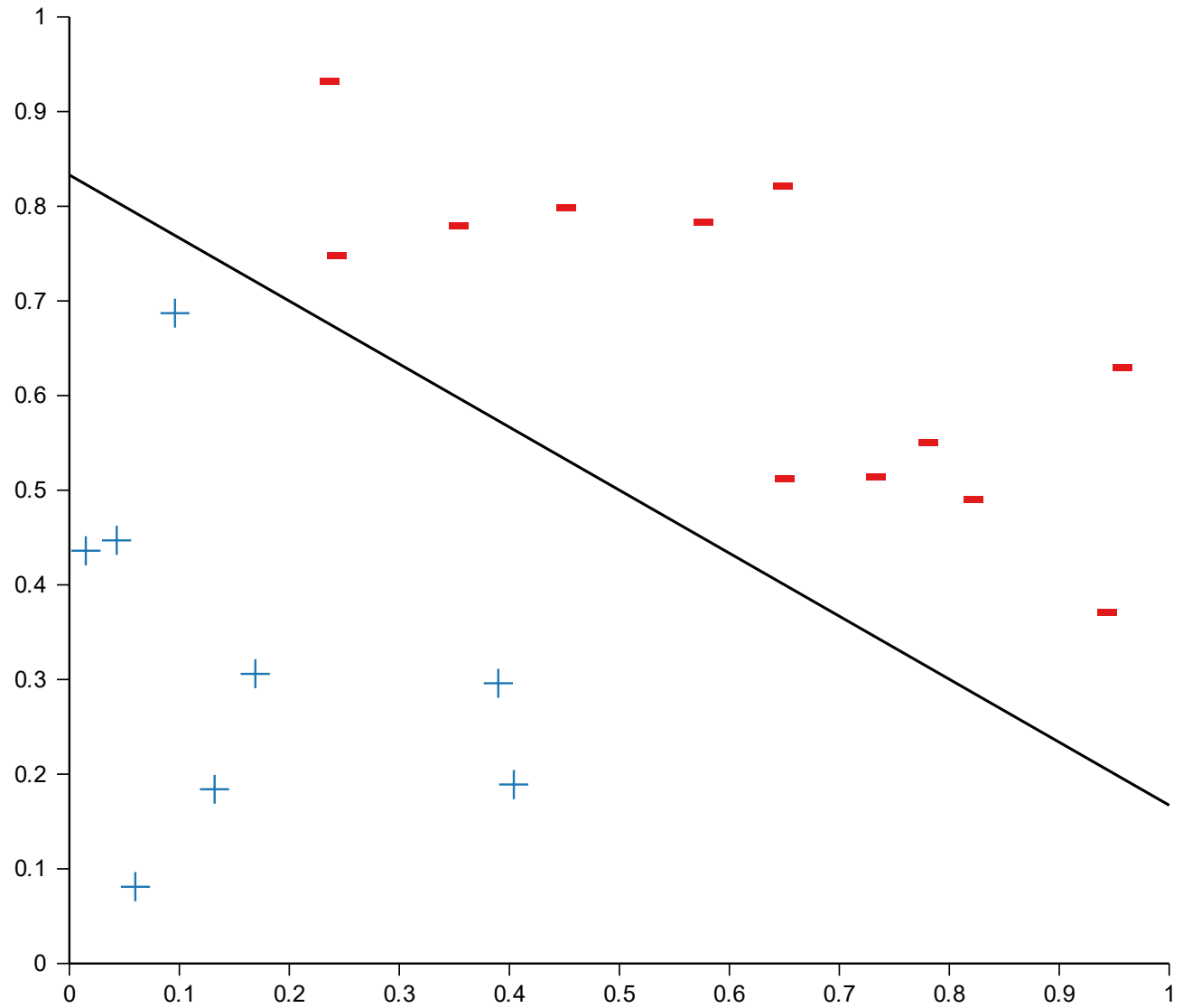
d. Increment t : $t \leftarrow t + 1$

• Output: $\mathbf{w}^{(t)}$

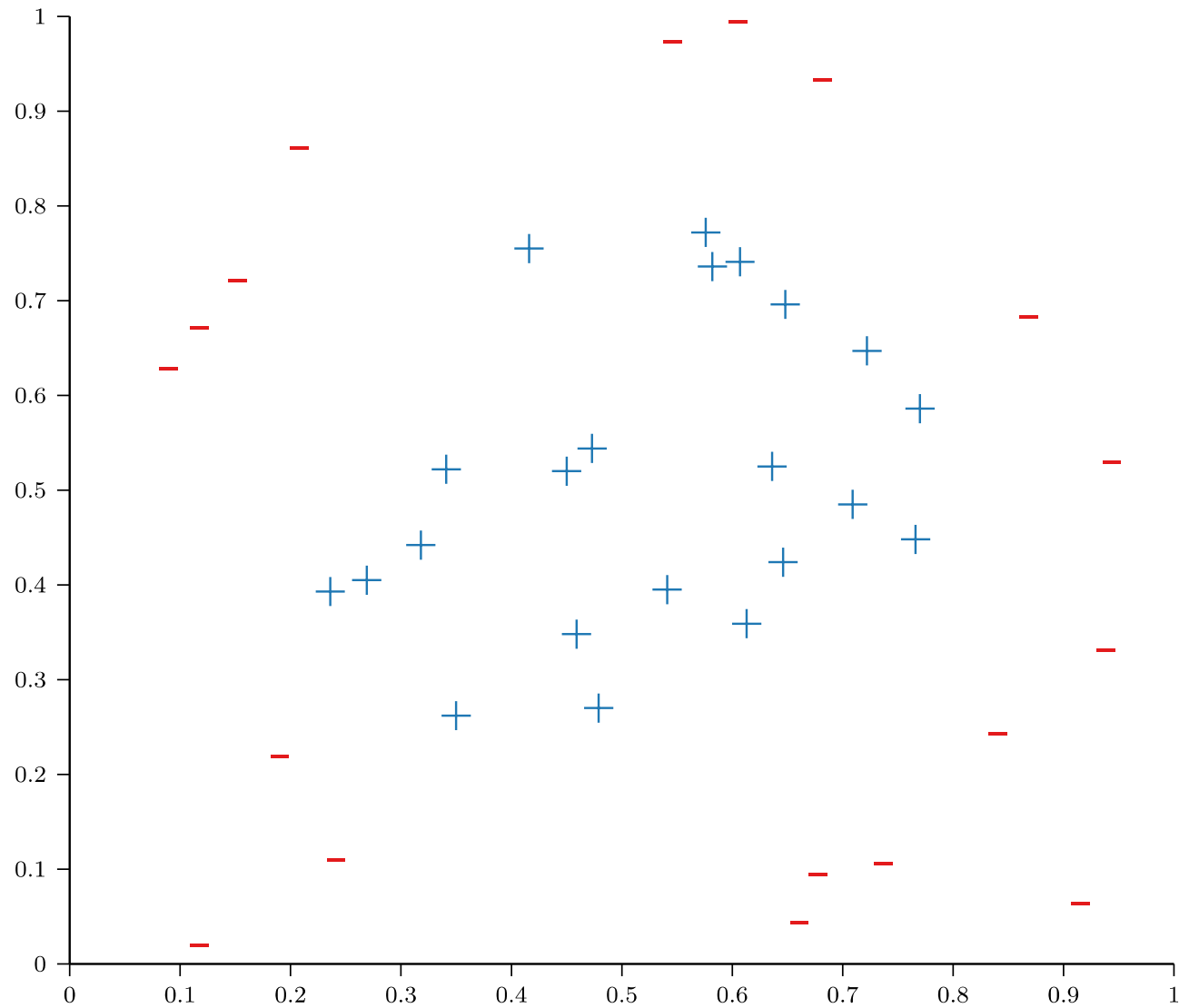
Linear Models



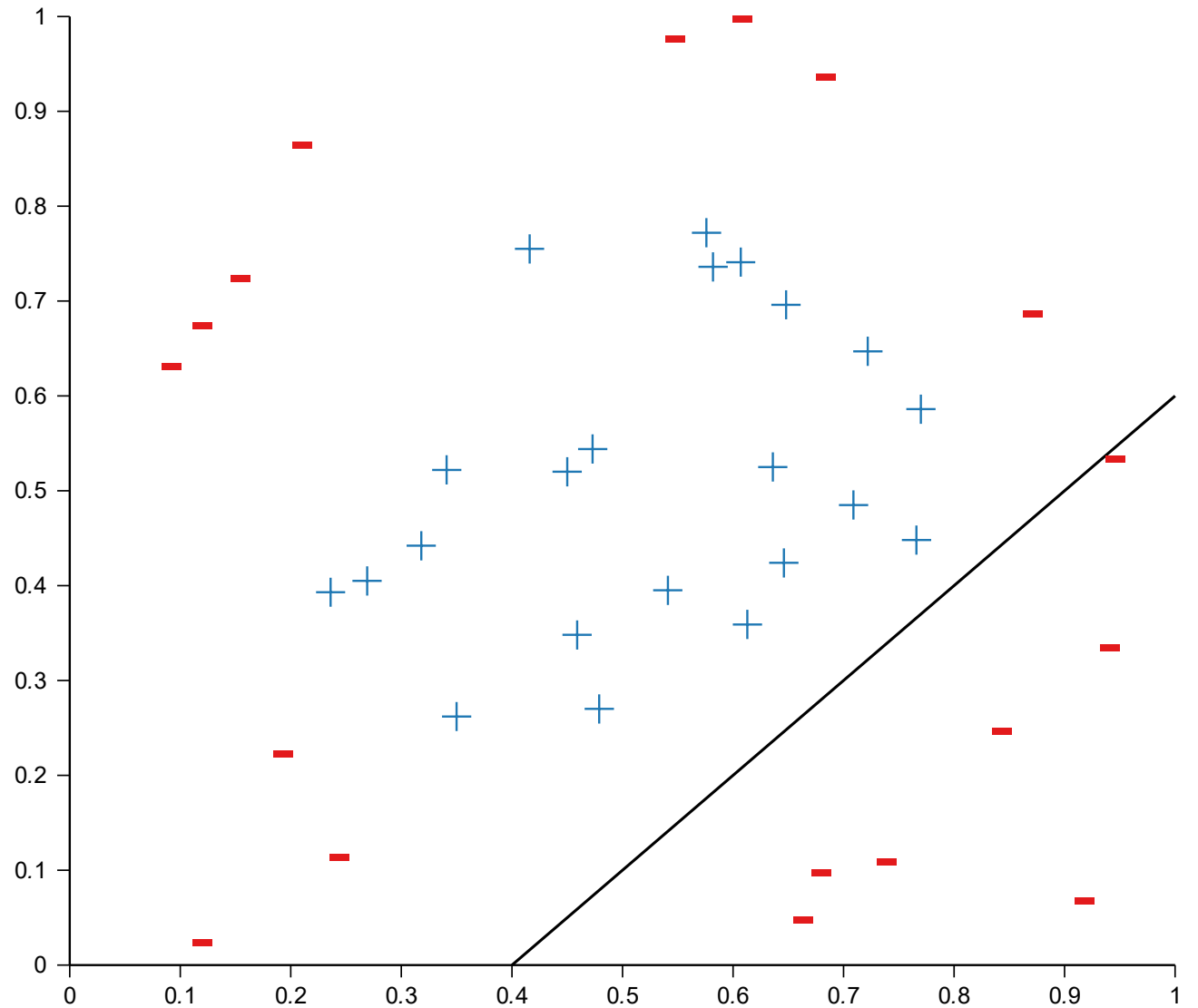
Linear Models



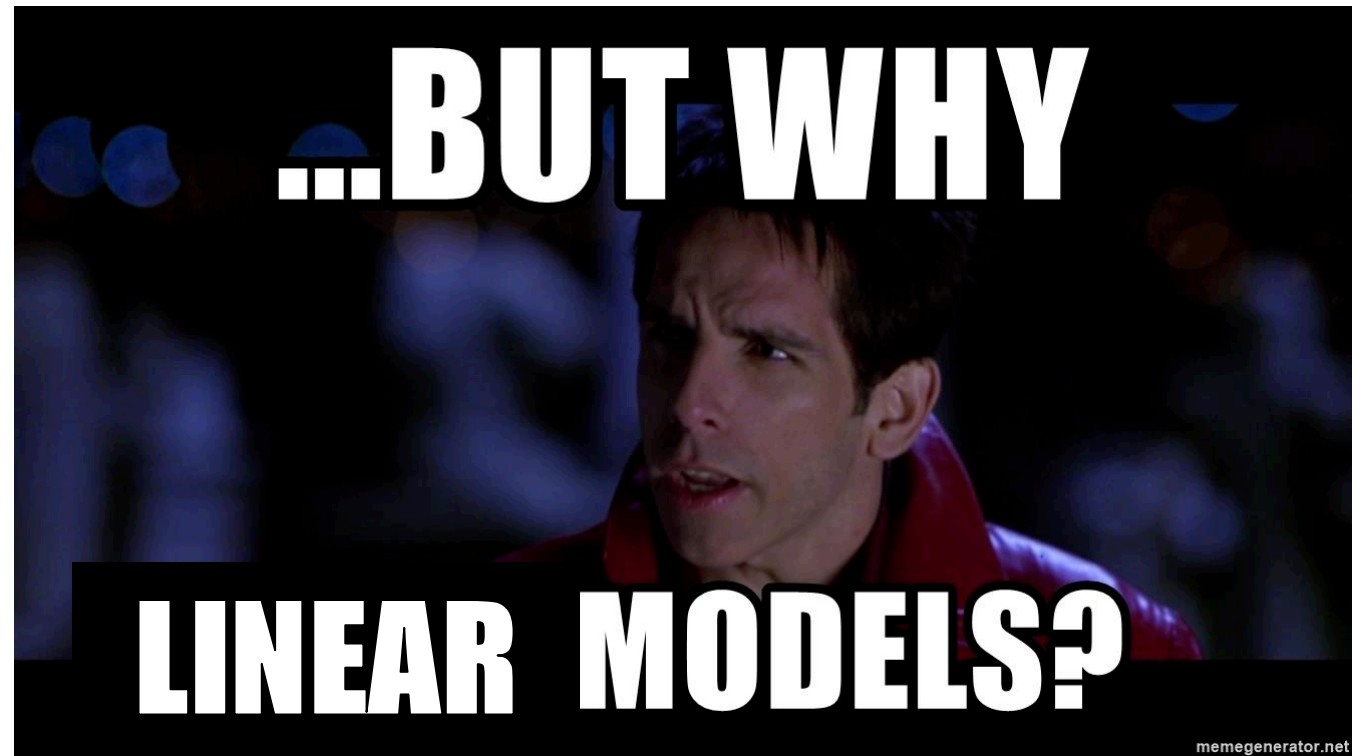
Linear Models?



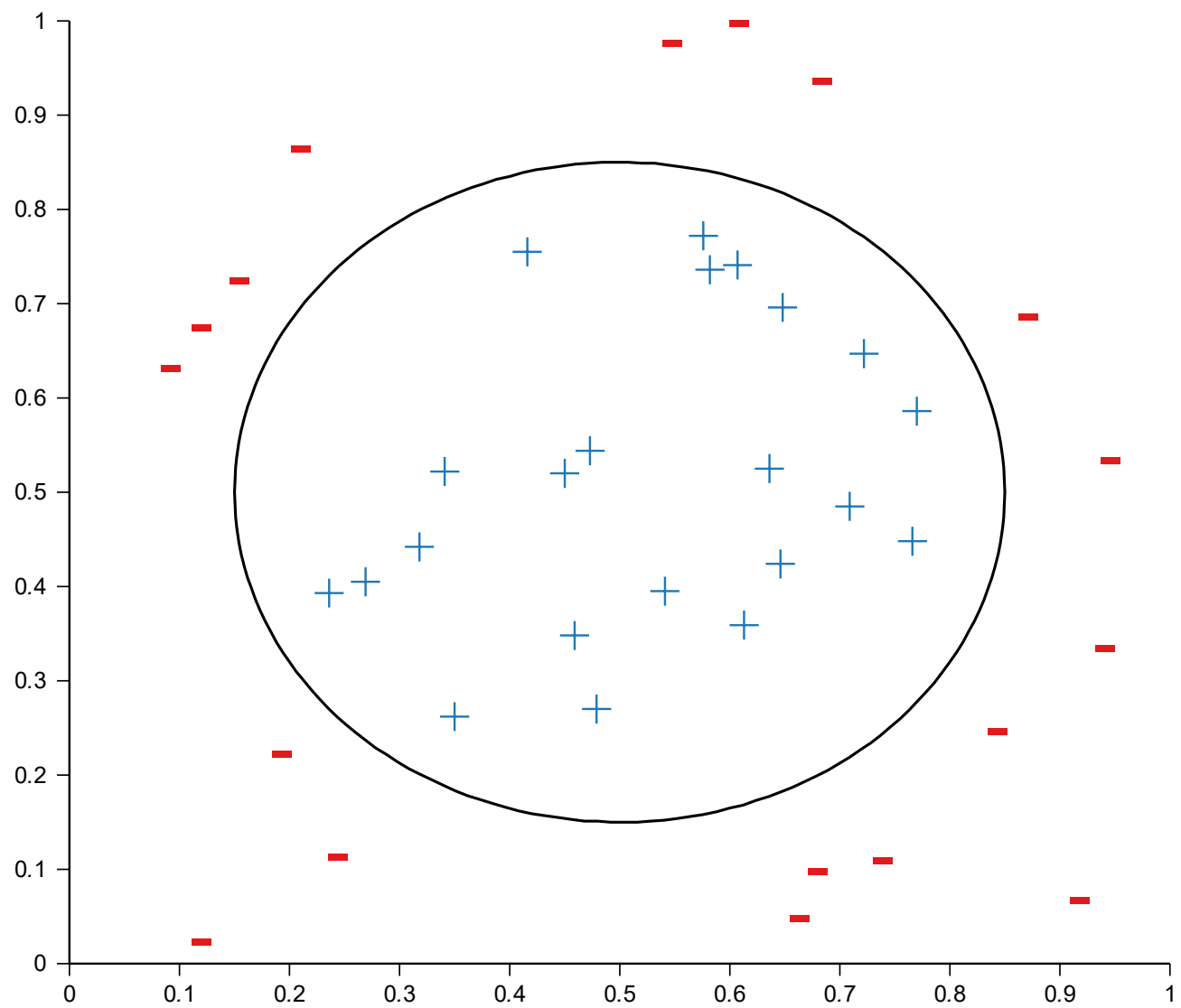
Linear Models?



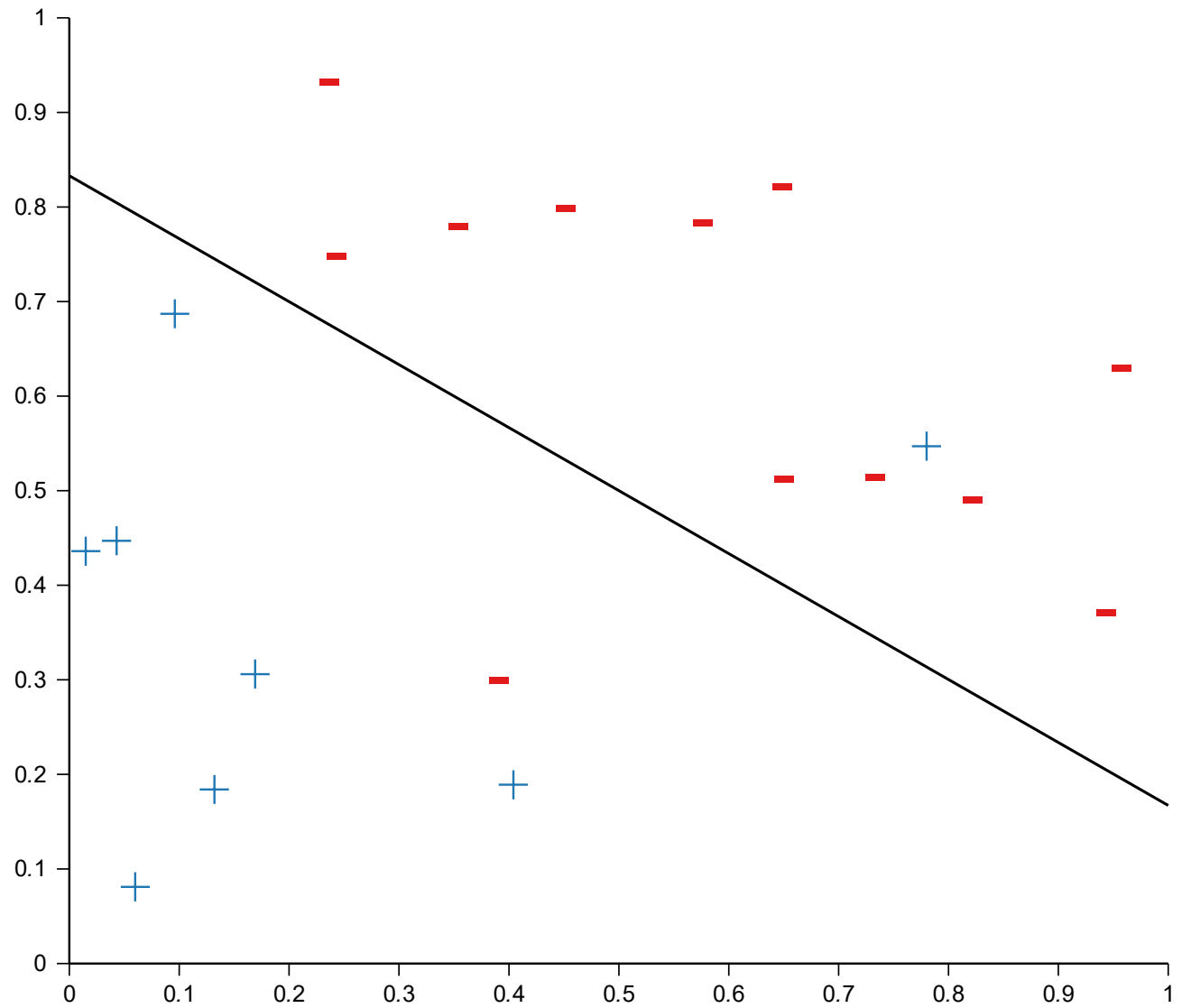
Linear Models?



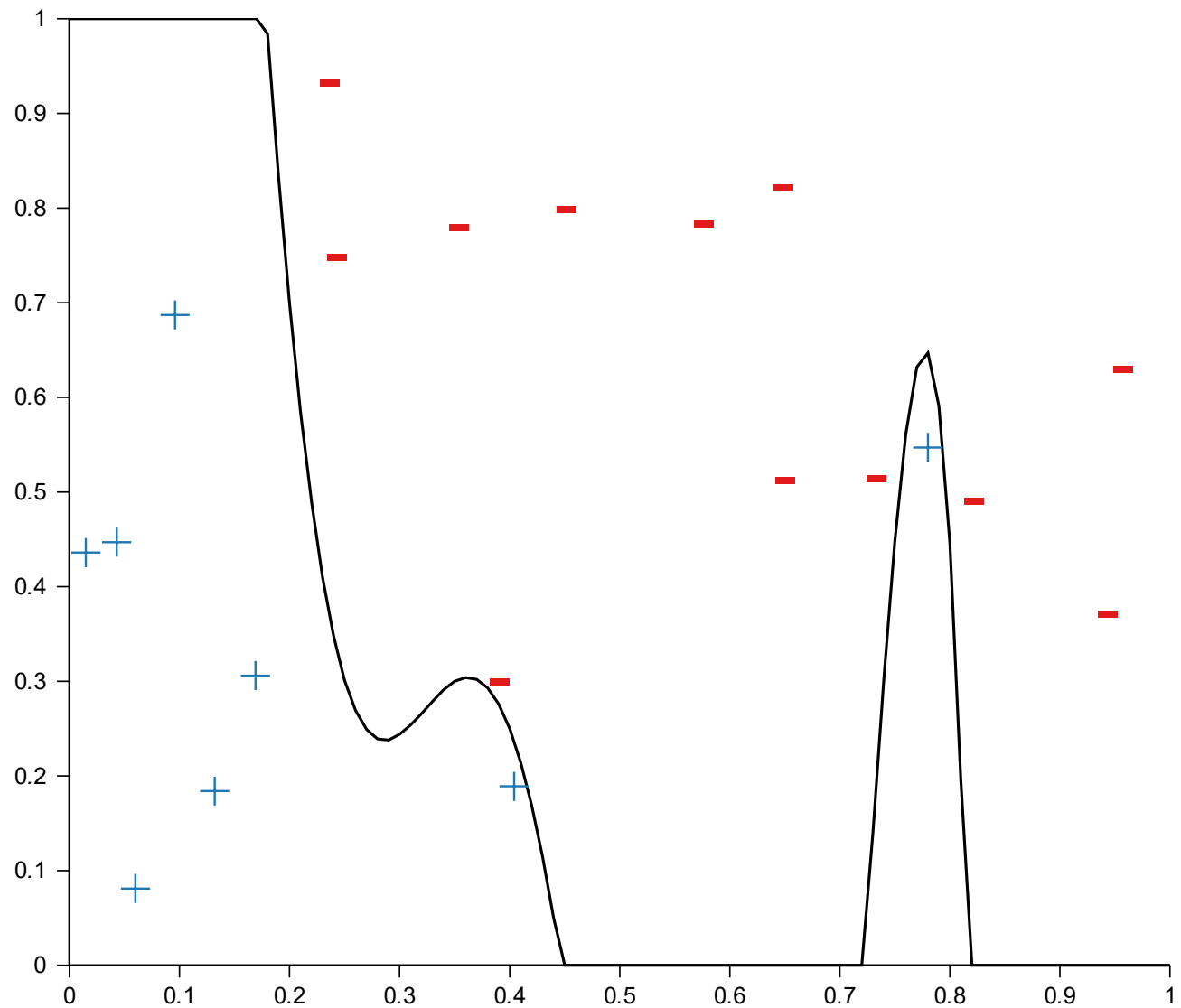
Nonlinear Models



Linear Models



Nonlinear Models?



Feature Transforms

- Given D -dimensional inputs $\mathbf{x} = [x_1, \dots, x_D]$, first compute some transformation of our input, e.g.,

$$\phi([x_1, x_2]) = [z_1 = (x_1 - 0.5)^2, z_2 = (x_2 - 0.5)^2]$$

General Q^{th} -order Transforms

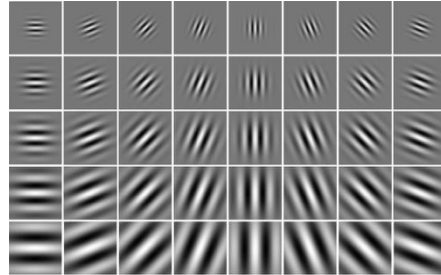
- $\phi_{2,2}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2]$
- $\phi_{2,3}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3]$
- $\phi_{2,4}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3, x_1^4, x_1^3x_2, x_1^2x_2^2, x_1x_2^3, x_2^4]$
- $\phi_{2,Q}$ maps a 2-D input to a $O(Q^2)$ -D output
- Scales even worse for higher-dimensional inputs...

Deterministic Transforms

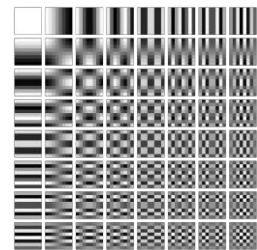
- log-transforms: $x \rightarrow \log(a + x)$
- power/root transforms: $x \rightarrow x^a$
- polynomial transforms: $x \rightarrow (x, x^2, x^3, \dots)$

Harmonic Transforms

- Gabor Filters:



- Discrete Cosine Transforms:



- ...

Can also have Linear Transforms

$$X \mapsto \{w_j^T X\}_{j=1}^k$$

- Principal Component Analysis
- Independent Component Analysis

Regression with basis functions

$$f(X) = \sum_{j=0}^m \beta_j \phi_j(X)$$

Basis coefficients

Basis functions (Linear combinations yield meaningful spaces)

Polynomial Basis

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

⋮

Fourier Basis

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

⋮

Wavelet Basis

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

⋮

Good representation for
periodic functions

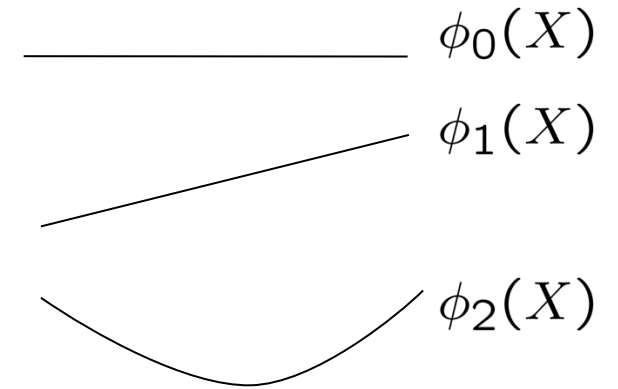
Good representation for
local functions

Regression with nonlinear features

$$f(X) = \sum_{j=0}^m \beta_j X^j = \sum_{j=0}^m \beta_j \phi_j(X)$$

Weight of
each feature

Nonlinear
features



In general, use any nonlinear features

e.g. e^X , $\log X$, $1/X$, $\sin(X)$, ...

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$$

or

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$\mathbf{A} = \begin{bmatrix} \phi_0(X_1) & \phi_1(X_1) & \dots & \phi_m(X_1) \\ \vdots & & \ddots & \vdots \\ \phi_0(X_n) & \phi_1(X_n) & \dots & \phi_m(X_n) \end{bmatrix}$$

$$\hat{f}_n(X) = \mathbf{X} \hat{\beta}$$

$$\mathbf{X} = [\phi_0(X) \ \phi_1(X) \ \dots \ \phi_m(X)]$$

Polynomial Regression

Univariate (1-dim)
case:

$$f(X) = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_m X^m = \mathbf{X}\beta$$

degree m
↙

$$\text{where } \mathbf{X} = [1 \ X \ X^2 \ \dots \ X^m] \quad \beta = [\beta_1 \ \dots \ \beta_m]^T$$

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \quad \text{or} \quad (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y} \quad \hat{f}_n(X) = \mathbf{X} \hat{\beta}$$

where

$$\mathbf{A} = \begin{bmatrix} 1 & X_1 & X_1^2 & \dots & X_1^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_n & X_n^2 & \dots & X_n^m \end{bmatrix}$$

Multivariate (p-dim)
case:

$$\begin{aligned} f(X) &= \beta_0 + \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)} \\ &+ \sum_{i=1}^p \sum_{j=1}^p \beta_{ij} X^{(i)} X^{(j)} + \sum_{i=1}^p \sum_{j=1}^p \sum_{k=1}^p X^{(i)} X^{(j)} X^{(k)} \\ &+ \dots \text{terms up to degree } m \end{aligned}$$

Mythbusters

Mythbusters

MYTH

If the data isn't linearly separable, just keep increasing the polynomial degree.

Mythbusters

Mythbusters

MYTH

If the data isn't linearly separable, just keep increasing the polynomial degree.

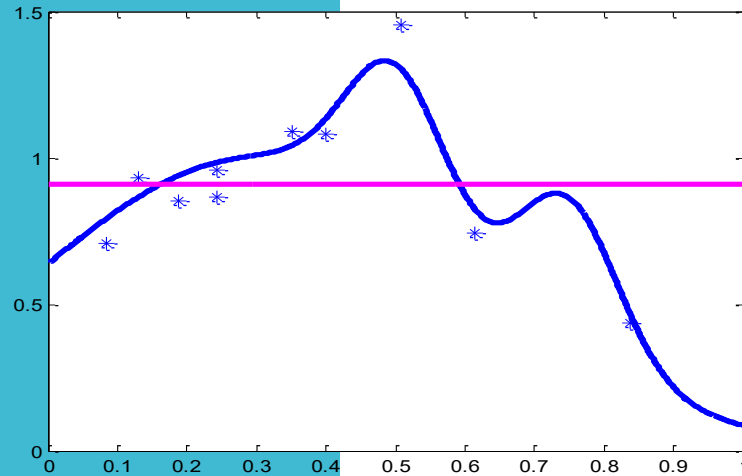
BUSTED

Higher-order feature transforms can overfit (bias-variance tradeoff) and can be very expensive to compute (e.g., $\Phi^{10}(x)$ costs $O(D^{10})$).

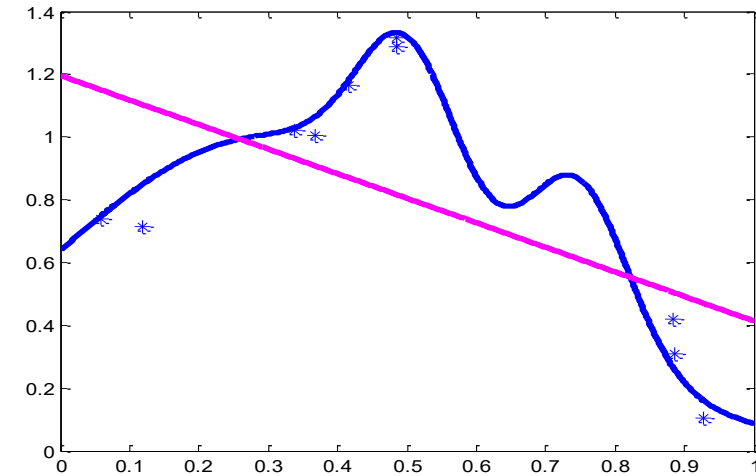
Polynomial Regression

Polynomial of order k , equivalently of degree up to $k-1$

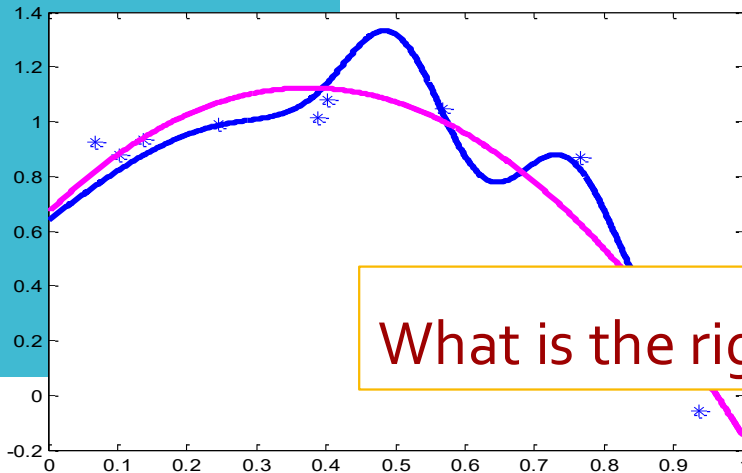
$k=1$



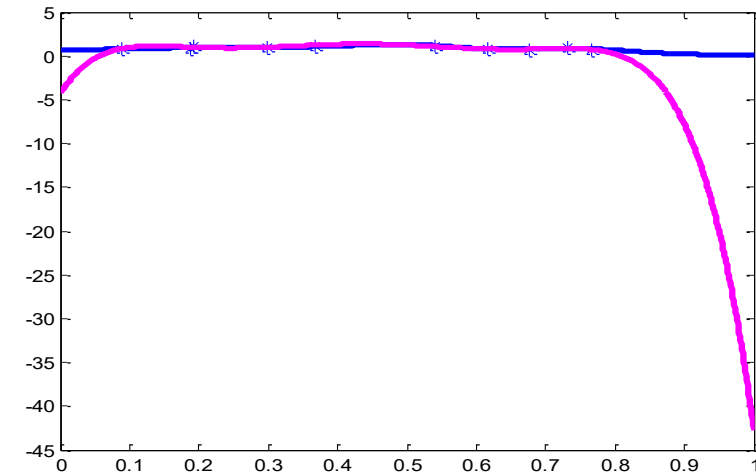
$k=2$



$k=3$

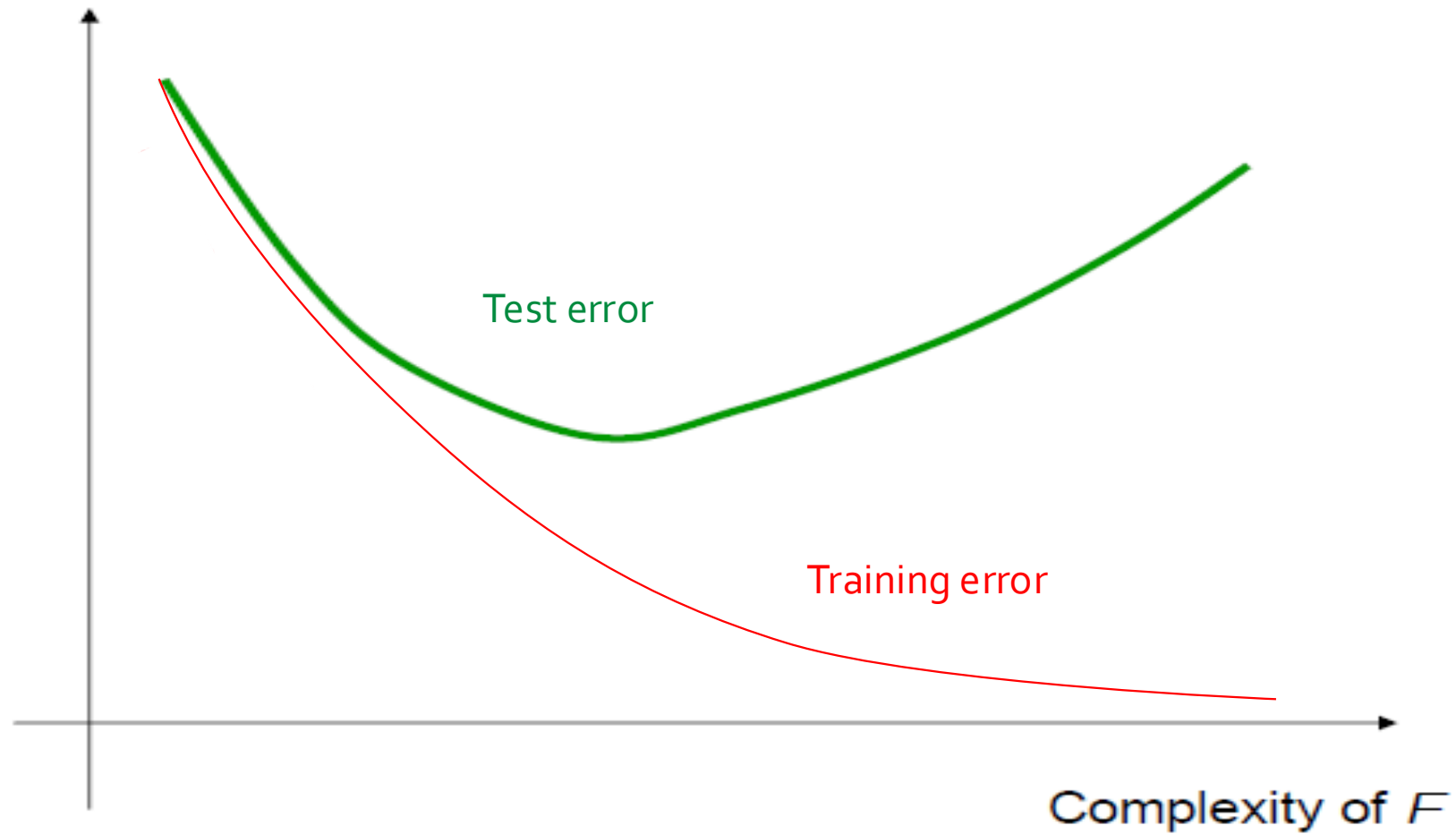


$k=7$



What is the right order?

Effect of Model Complexity



Efficiency

- Depending on the transformation Φ and the dimensionality of the original input space \mathcal{X} , computing $\Phi(\mathbf{x})$ can be prohibitively computationally expensive
 - Computing $\Phi_2(\mathbf{x}) = [x_1, x_2, \dots, x_D, x_1^2, x_1x_2, \dots, x_D^2]$ for $\mathbf{x} \in \mathbb{R}^D$ requires $D + \binom{D}{2} + D = \frac{D^2+3D}{2} = O(D^2)$ time
 - Computing $\Phi_{10}(\mathbf{x})$ requires $O(D^{10})$ time
- Tradeoff:
 - High-dimensional transformations can result in good hypotheses (as long as they don't overfit) but...
 - High-dimensional transformations are expensive