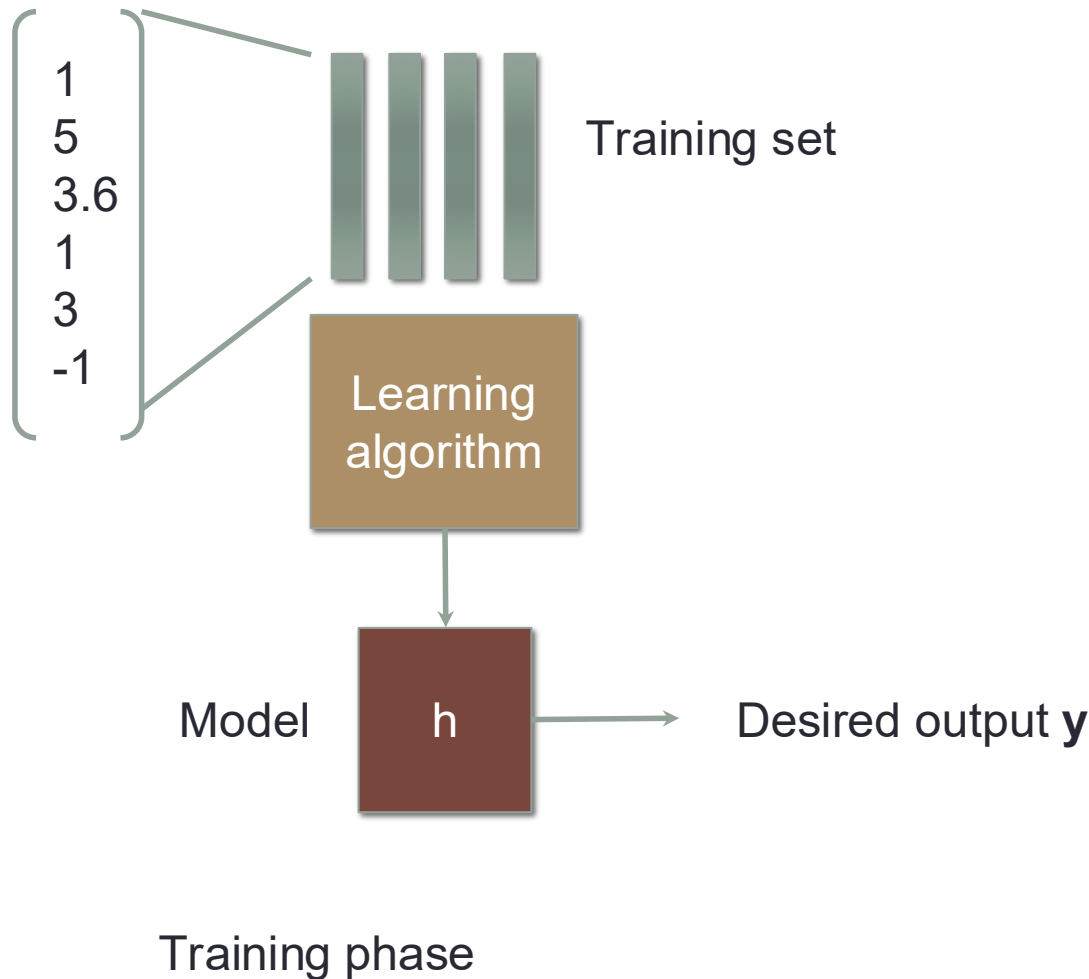


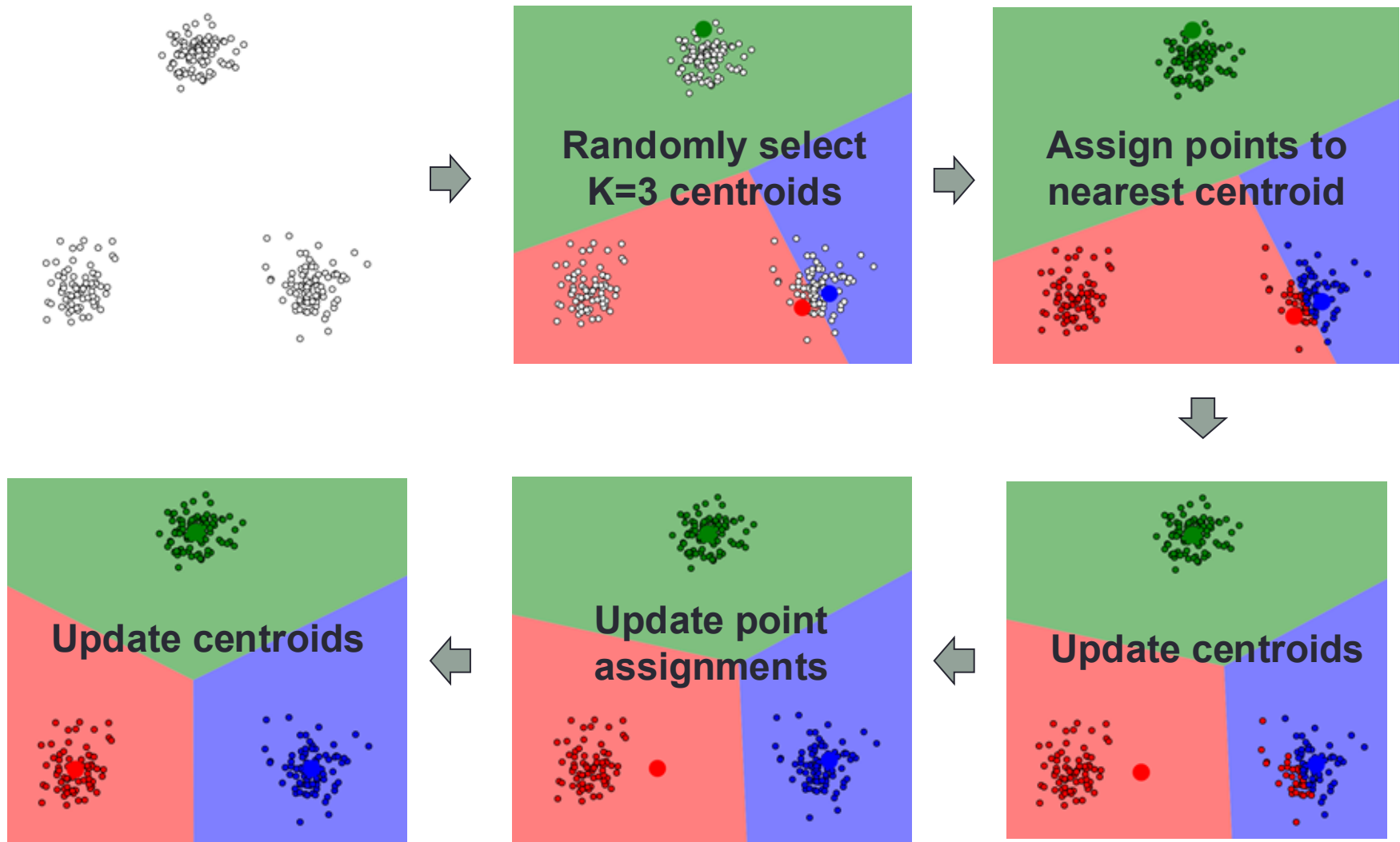
REGRESSION

with some MLE

How do we learn from data?



An Illustration Of K-Mean Clustering



Types of machine learning

1. Supervised learning

Learn a model F from pairs of (\mathbf{x}, y)

2. Unsupervised learning

Discover the hidden structure in unlabeled data \mathbf{x} (no y)

3. Reinforcement learning

Train an agent to take appropriate actions in an environment by maximizing rewards

REGRESSION

(Linear) Regression

- $h_{\theta}(\mathbf{x}_1) = \theta_0 + \theta_1 x_{1,1} + \theta_2 x_{1,2} + \theta_3 x_{1,3} + \theta_4 x_{1,4} + \theta_5 x_{1,5}$

- θ s are the parameter (or weights)

Assume x_0 is always 1

- We can rewrite n is dimension of x

$$h_{\theta}(\mathbf{x}_i) = \sum_{j=0}^n \theta_j x_{i,j} = \theta^T \mathbf{x}_i$$

h is parameterized by θ

- Notation: vectors are bolded
- Notation: vectors are column vectors

Cost function (Loss function)

- Let's use the mean square error (MSE)

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i)^2$$

m is the number of training examples



We want to pick θ that minimize the loss

i here is the index of the training example

Note how \mathbf{x} is bolded

LMS regression with gradient descent

$$\frac{\partial J}{\partial \theta_j} = -\sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i) x_i^{(j)}$$

$$\theta_j \Leftarrow \theta_j + r \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i) x_i^{(j)}$$

Or using the close-form solution

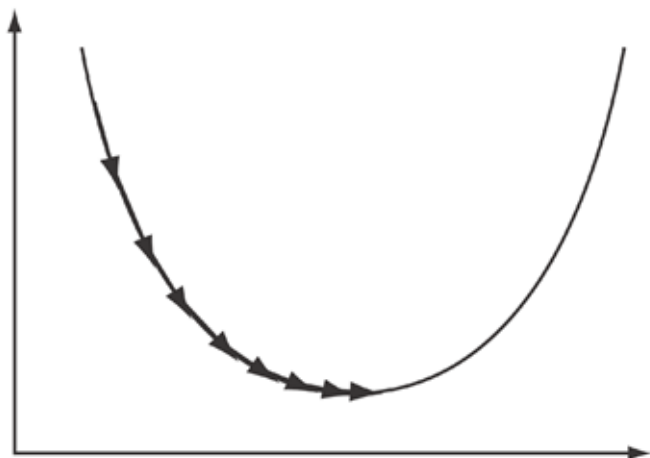
$$\theta = (X^T X)^{-1} X^T y$$

Loose Ends from HW

- How to select r ? (The learning rate)

$$\theta_j \Leftarrow \theta_j + r \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i) x_i^{(j)}$$

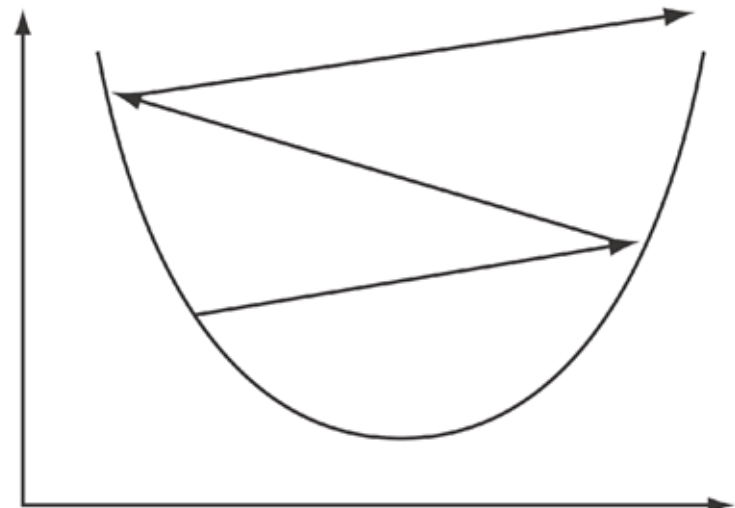
Small learning rate



<https://www.packtpub.com/books/content/big-data>

r too small and the model converges slowly

Large learning rate



r too large and the model diverges

Learning rate issues

- Typically, r is normalized with the amount of training examples in a mini-batch. (Divide by m)

$$\theta_j \Leftarrow \theta_j + r \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i) x_i^{(j)}$$

- Typical values are 0.1-0.001
- Usually have a decay over time

Scaling the input data

- We use age, passenger class, gender, and embark as our input.
- Age has a lot more variance (0.42 – 80) than the other data.
- This makes parameter initialization hard and makes the learning rate selection hard.
- $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

Scaling the input data

- Scale all input data to be in the same range
- Using statistics from training data
 - Scale to $[-1, 1]$
 - Scale to $[0, 1]$
 - Scale to standard normal
- Don't forget to apply the same scaling to the test data

Feature selection

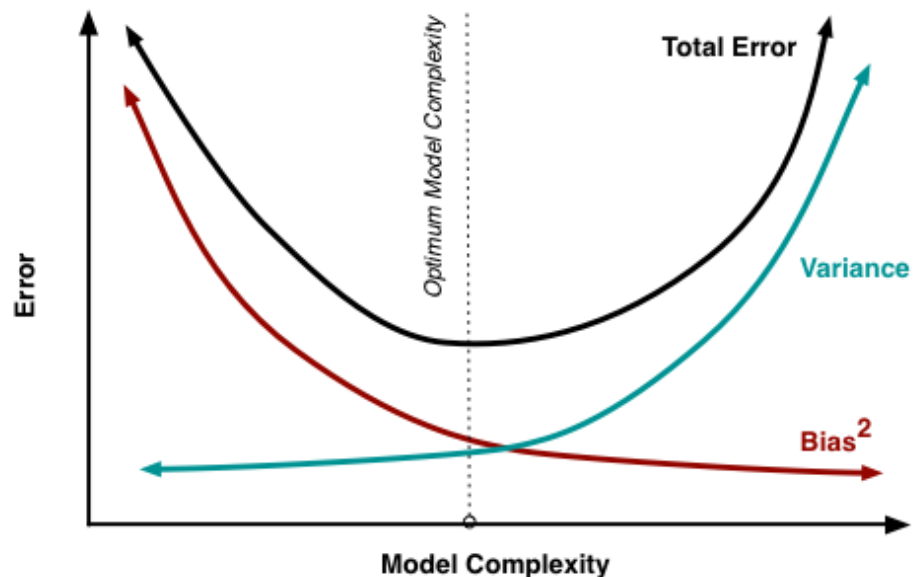
- Most likely you will get better results with just two features.
- This is the importance of feature selection.
- Knowing what good features to select is not trivial
- Approaches for feature selection (or for not having to do feature selection)
 - Cross validation
 - Random forest
 - Boosting

When to stop the update?

- Consider the updates of Logistic regression as trying to reduce the bias of the model
 - As we keep updating, the model overfits more to the training data
- We want to stop when the error on the validation set increases*
- More on this later


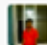



- Validation test: a separate set that is used to measure overfitting

Training set
Validation set
Test set



More tricks?

- <http://ahmedbesbes.com/how-to-score-08134-in-titanic-kaggle-challenge.html>
- Feature Engineering/selection
- Parameter tuning
- Try different models

449	▲ 62...	Kaustubh Kulkarni 2		0.81340	6	6h
450	new	AshishDoshi		0.81340	1	5h
451	new	SouravKarwa		0.81340	2	31m
452	▲ 18...	Ahmed Besbes		0.81340	15	now
Your Best Entry ↑ Your submission scored 0.81340, which is not an improvement of your best score. Keep trying!						
453	▼ 7	Clement Sengelen		0.80861	11	2mo

Other loss functions

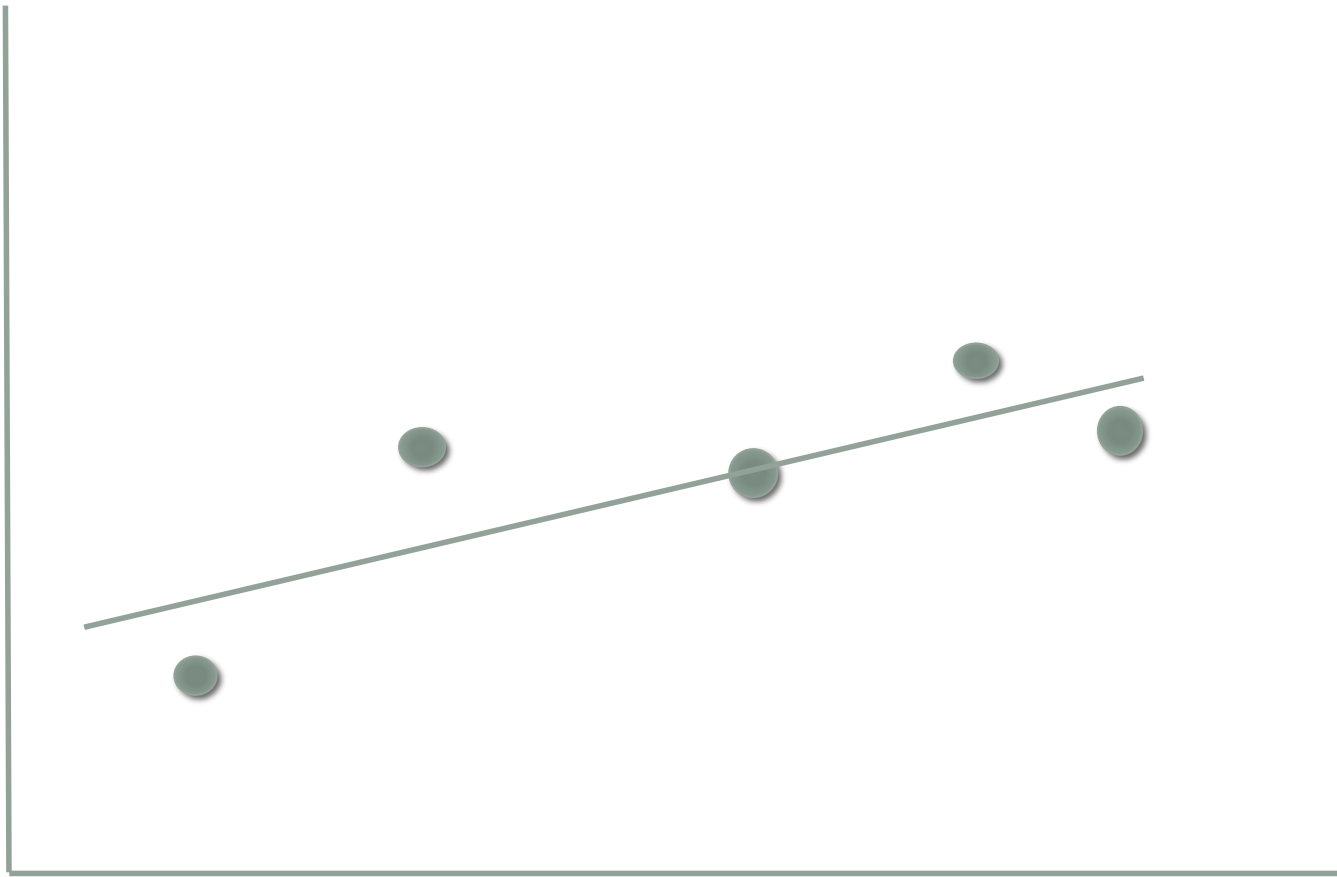
- MSE

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i)^2$$

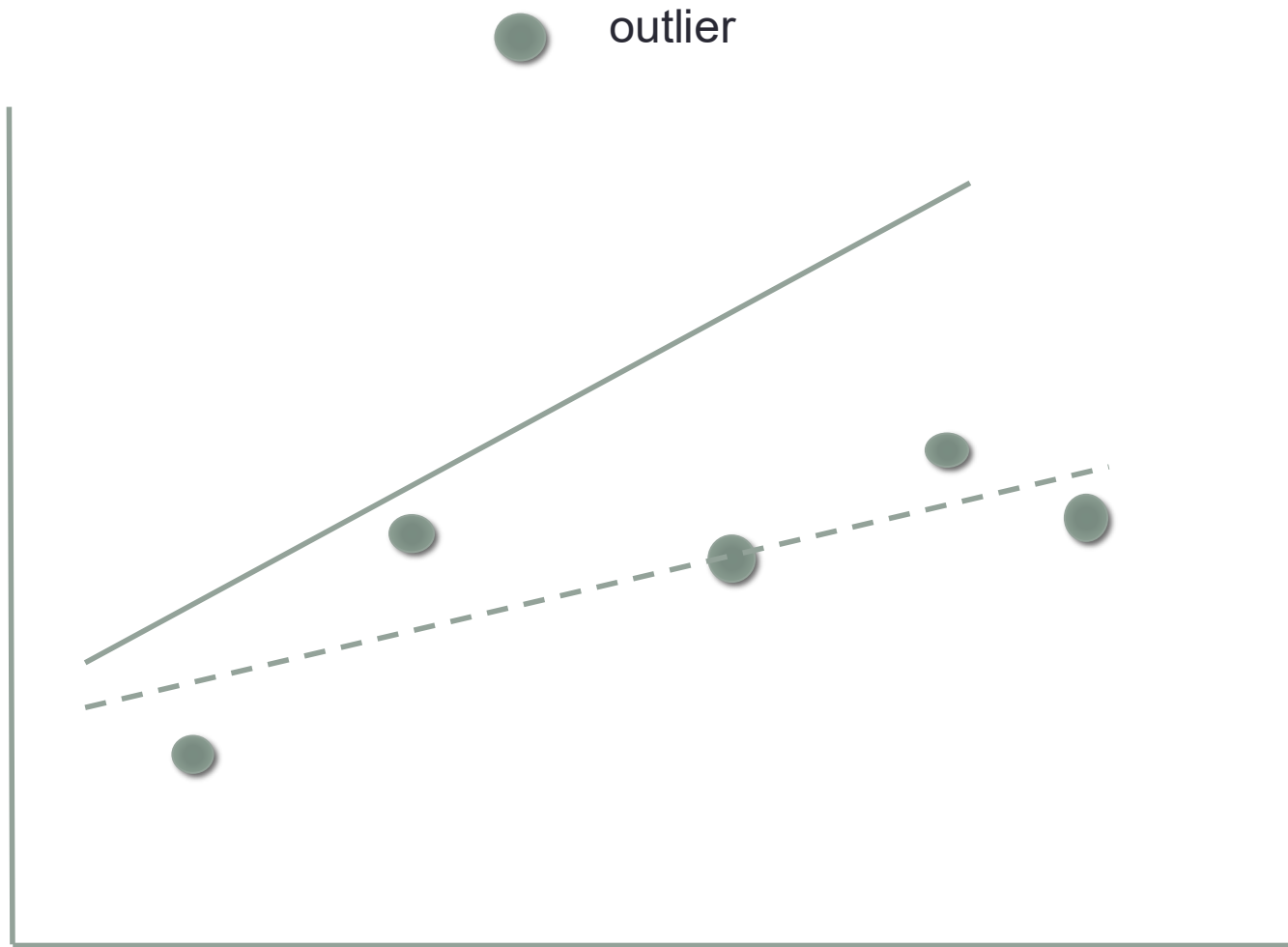
- Also called L2 loss
- L1 loss

$$\frac{1}{m} \sum_{i=1}^m |y_i - \theta^T \mathbf{x}_i|$$

L2 vs L1 loss



L2 vs L1 loss



Outlier frequently happens in the real world

Norms (p-norm or Lp-norm)

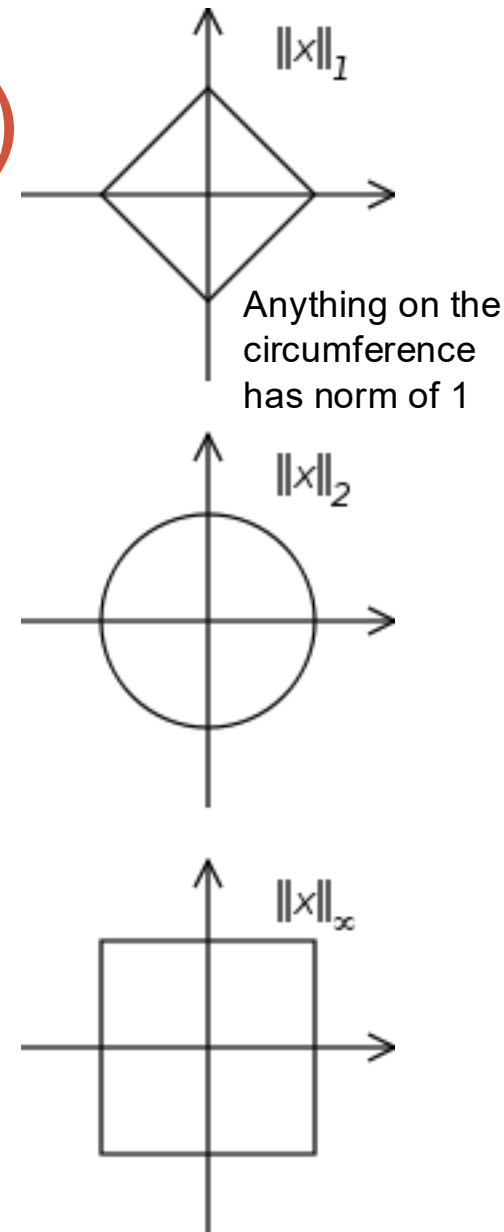
- For any real number $p > 1$

$$\|\mathbf{x}\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}}$$

- For $p = \infty$

$$\|x\|_{\infty} = \max \{|x_1|, |x_2|, \dots, |x_n|\}$$

- We'll see more of p-norms when we get to neural networks



Predicting floods

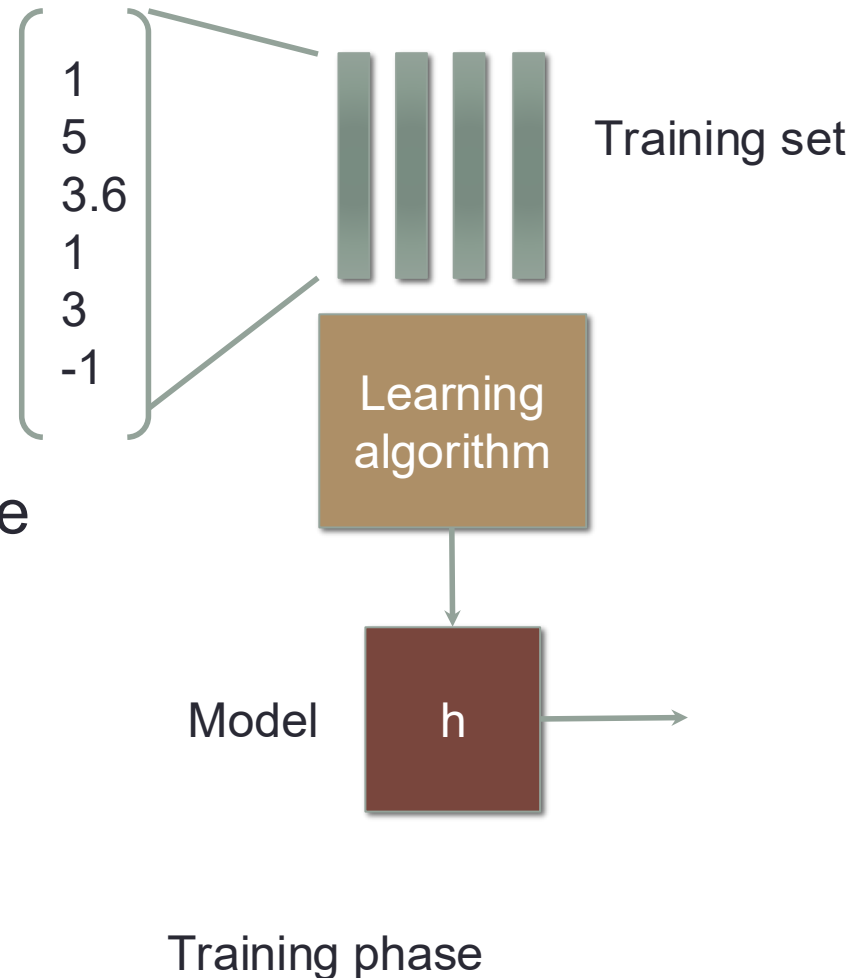
Cloth	Corn	Grass	Water	Beer	Flood?
4	6	3	10	0	yes
5	1	0	0	7	yes
6	0	3	5	7	no
5	0	3	12	0	yes
4	3	0	6	7	?

So far we talk about predicting an amount what if we want to do classification

Let's start with a binary choice. Flood or no flood

Flood or no flood

- What would be the output?
- $y = 0$ if not flooded
- $y = 1$ if flooded
- Anything in between is a score for how likely it is to flood

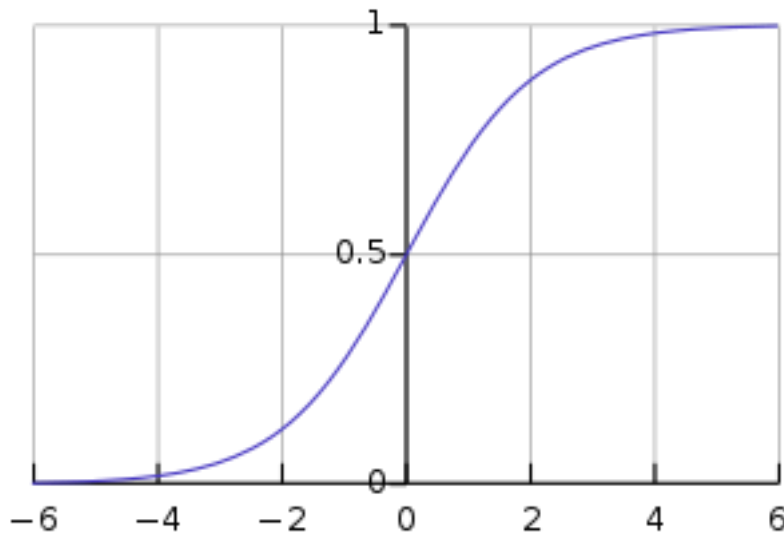


Can we use regression?

- Yes
- $h_{\theta}(\mathbf{x}_1) = \theta_0 + \theta_1 x_{1,1} + \theta_2 x_{1,2} + \theta_3 x_{1,3} + \theta_4 x_{1,4} + \theta_5 x_{1,5}$
- But
- What does it mean when h is higher than 1?
- Can h be negative? What does it mean to have a negative flood value?

Logistic function

- Let's force h to be between 0 and 1 somehow
- Introducing the logistic function (sigmoid function)



$$\begin{aligned} f(x) &= \frac{1}{1 + e^{-x}} \\ &= \frac{e^x}{1 + e^x} \end{aligned}$$

Logistic Regression

- Pass $\theta^T \mathbf{x}$ through the logistic function

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Loss function?

- MSE error no longer a good candidate

Distribution parameter estimation

- $P(\text{head}) = \theta$, $\theta = \text{\#heads}/\text{\#tosses}$
- HHTTH

-

- $L(\theta) = P(X; \theta) = P(\text{HHTTH}; \theta)$
- Maximum Likelihood Estimate (MLE)
 - Likelihood - Probability of encountering the data X given the parameters θ

Probabilistic Interpretation of linear regression

- Real world data is our model plus some error term
 - Noise in the data
 - Something that we do not model (features we are missing)
- Let's assume the error is normally distributed with mean zero and variance σ^2
 - Why Gaussian?
 - Why saying mean is zero is a valid assumption?

$$y_i = \theta^T \mathbf{x}_i + \epsilon_i$$

Probabilistic view of Linear regression

From our assumption we know that

$$y_i = \theta^T \mathbf{x}_i + \epsilon_i$$

$$p(y_i | \mathbf{x}_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\overbrace{(y_i - \theta^T \mathbf{x}_i)^2}^{\text{error}}}{2\sigma^2}\right)$$

Error term is normally distributed with mean 0 and variance σ^2

We want to find θ

Maximize Likelihood of seeing x and y in training

or in other words maximize the probability of
answering y from x in the training data

$$p(y_i | \mathbf{x}_i; \theta)$$

Maximizing Likelihood

What is the assumption here?
Is it accurate?

- Max $L(\theta) = \prod_{i=1}^m p(y_i | \mathbf{x}_i; \theta)$
We use the log likelihood instead $\log(L(\theta)) = l(\theta)$

From our previous lecture

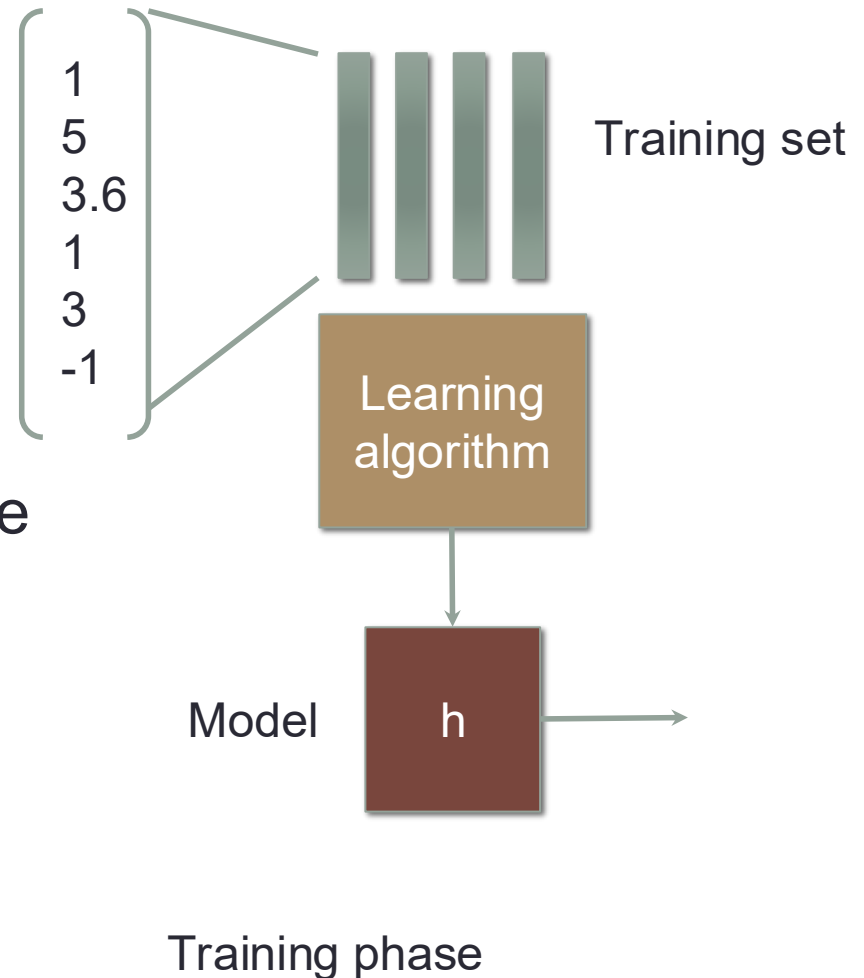
$$\text{Min } J(\theta) = \frac{1}{m} \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i)^2$$

Mean square error solution and MLE solution

- Turns out MLE and MSE gets to the same solution
 - This justifies our choice of MSE as the Loss for linear regression
 - This does not mean MSE is the best Loss for regression, but you can at least justify it under this probabilistic reasoning and assumptions
- Note how our choice of variance σ^2 falls out of the maximization, so this derivation is true regardless of which size of the variance.
- Note that MLE derivation assumes that the error is normally distributed! **This is a key assumption for linear regression.**
 - Error is normally distributed is not that same as y is normally distributed.

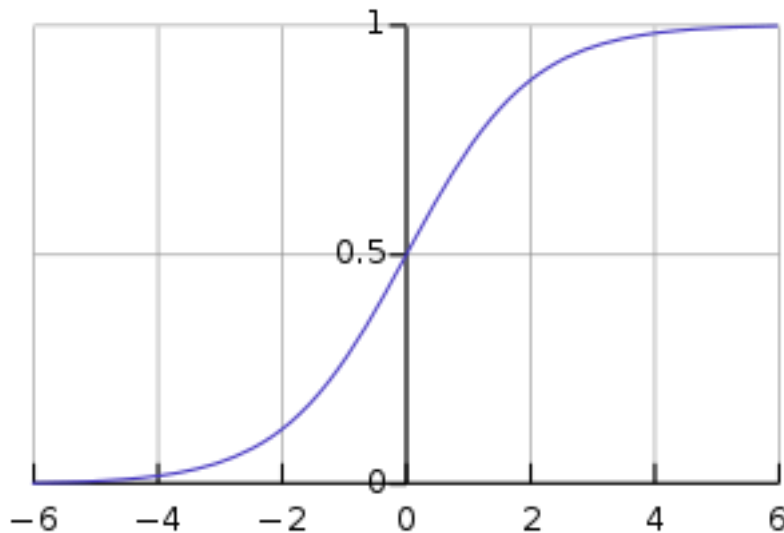
Flood or no flood

- What would be the output?
- $y = 0$ if not flooded
- $y = 1$ if flooded
- Anything in between is a score for how likely it is to flood



Logistic function

- Let's force h to be between 0 and 1 somehow
- Introducing the logistic function (sigmoid function)



$$\begin{aligned} f(x) &= \frac{1}{1 + e^{-x}} \\ &= \frac{e^x}{1 + e^x} \end{aligned}$$

Logistic Regression

- Pass $\theta^T \mathbf{x}$ through the logistic function

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Loss function?

- MSE error no longer a good candidate
- Let's turn to use probabilistic argument for logistic regression

Logistic Function derivative

The derivative has a nice property by design.

This is also why many algorithm we'll learn later in class also uses the logistic function

$$\begin{aligned} g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\ &= \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\ &= \frac{1}{(1 + e^{-z})} \cdot \left(1 - \frac{1}{(1 + e^{-z})} \right) \\ &= g(z)(1 - g(z)). \end{aligned}$$

Probabilistic view of Logistic Regression

- Let's assume, we'll classify as 1 with probability according to the output of

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$P(y = 1 \mid x; \theta) = h_{\theta}(x)$$

$$P(y = 0 \mid x; \theta) = 1 - h_{\theta}(x)$$

or

$$p(y \mid x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

Maximizing log likelihood

$$p(y \mid x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

$$\begin{aligned} g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\ &= \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\ &= \frac{1}{(1 + e^{-z})} \cdot \left(1 - \frac{1}{(1 + e^{-z})}\right) \\ &= g(z)(1 - g(z)). \end{aligned}$$

Logistic Regression update rule

$$\theta_j \Leftarrow \theta_j + r \sum_{i=1}^m (y_i - h_{\theta}(x_i)) x_i^{(j)}$$

Update rule for linear regression

$$\theta_j \Leftarrow \theta_j + r \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i) x_i^{(j)}$$

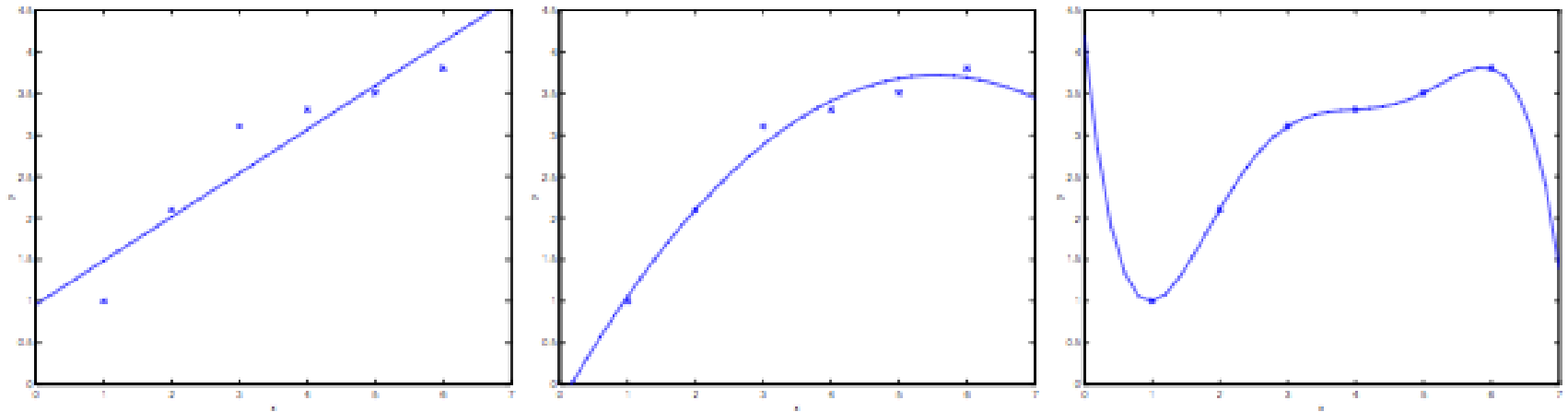
Regression with non-linear features

- If we add extra features that are non-linear
 - For example, x^2

Cloth	Corn	Grass	Water	Beer	Rainfall
4	6	3	10	0	76950
5	1	0	0	7	30234
6	0	3	5	7	123456
5	0	3	12	0	89301
4	3	0	6	7	?

- $h_{\theta}(\mathbf{x}_1) = \theta_0 + \theta_1 x_{1,1} + \theta_2 x_{1,2} + \theta_3 x_{1,3} + \theta_4 x_{1,4} + \theta_5 x_{1,5} + \theta_6 x_{1,1}^2 + \dots$
- These can be considered as additional features
- We can now have a line that is non-linear

Overfitting Underfitting



Adding more non-linear features makes the line more curvy
(Adding more features also means more model parameters)

The curve can go directly to the outliers with enough parameters.

We call this effect **overfitting**

For the opposite case, having not enough parameters to model the data is called **underfitting**

Feature engineering

- Logistic regression is a linear classification



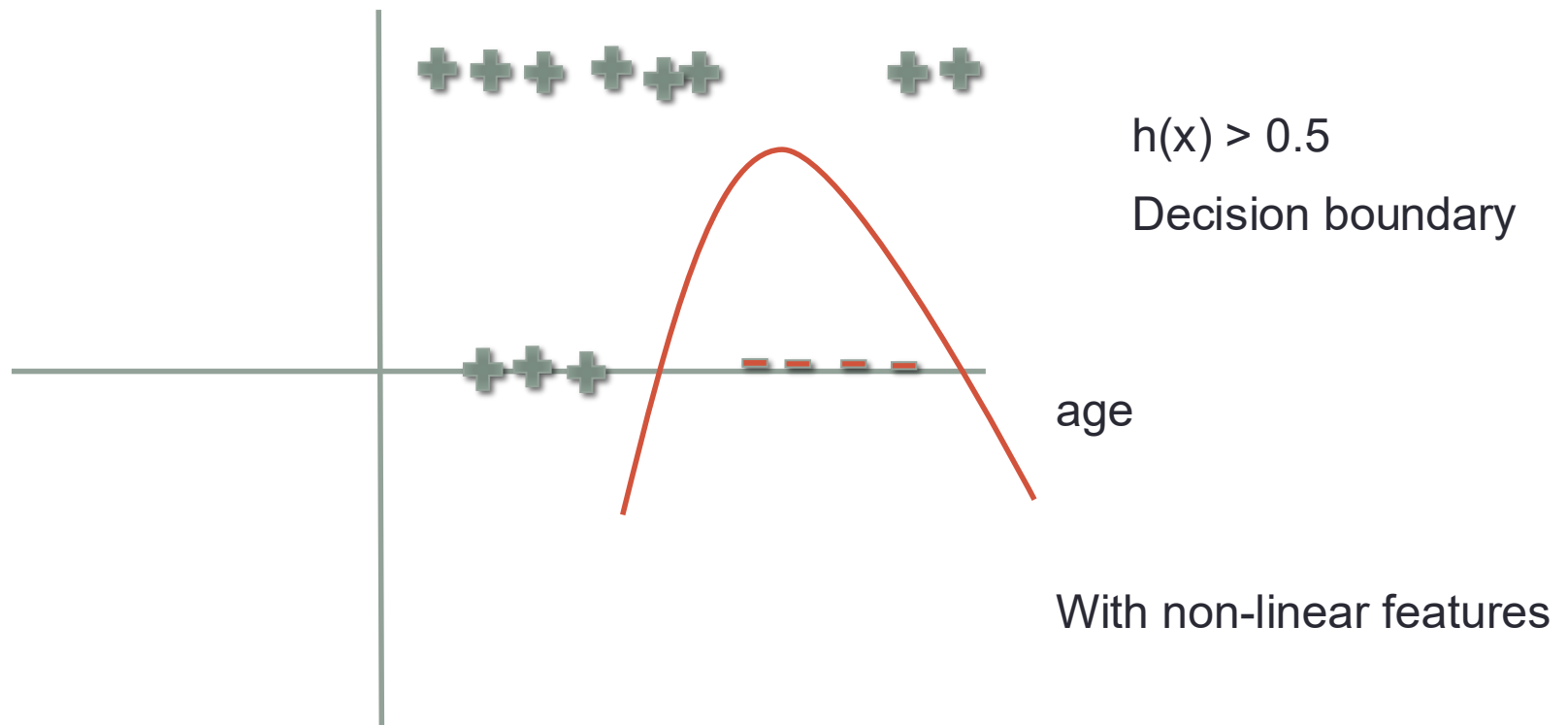
Feature engineering

- Logistic regression is a linear classification



Feature engineering

- Add non-linear features to get non-linear decision boundaries



This is also a form of feature selection (more specifically feature engineering)

Bias-Variance trade-off

- We will formulate overfitting and underfitting mathematically
- Using regression model

Regression with Gaussian noise

- $y = h(\mathbf{x}) + \varepsilon$
 - Where ε is normally distributed with mean zero and variance σ^2
 - The training data $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3) \dots\}$ is drawn from some distribution $P(\mathbf{x}, y)$ governing our universe!
 - Assume (\mathbf{x}_i, y_i) is iid
- Given D we can train a regressor $h_D(\mathbf{x})$
- We calculate the expected error (squared error) on new (\mathbf{x}, y) data with the regressor

$$\bullet E_{(\mathbf{x}, y)}[(h_D(\mathbf{x}) - y)^2] = \int \int_{\mathbf{x} \ y} (h_D(\mathbf{x}) - y)^2 \text{Pr}(\mathbf{x}, y) \partial y \partial \mathbf{x}$$

- But D is actually random too!

Regression with Gaussian noise

- We calculate the expected error (squared error) on new (\mathbf{x}, y) data with the regressor

- $$E_{(\mathbf{x}, y)}[(h_D(\mathbf{x}) - y)^2] = \int_{\mathbf{x}} \int_y (h_D(\mathbf{x}) - y)^2 \text{Pr}(\mathbf{x}, y) \partial y \partial \mathbf{x}$$

- Consider parallel worlds, we can receive different training data D which yields different regression $h_D(\mathbf{x})$
- The expectation of error over all possible new test data point (\mathbf{x}, y) and different possible training data D is

$$E_{\substack{(\mathbf{x}, y) \sim P \\ D \sim P^n}} [(h_D(\mathbf{x}) - y)^2] = \int_D \int_{\mathbf{x}} \int_y (h_D(\mathbf{x}) - y)^2 P(\mathbf{x}, y) P(D) \partial \mathbf{x} \partial y \partial D$$

Regression with Gaussian noise

- This expression tells the expected quality of our model with random training data and a random test data

$$E_{\substack{(\mathbf{x}, y) \sim P \\ D \sim P^n}} \left[(h_D(\mathbf{x}) - y)^2 \right] = \int_D \int_{\mathbf{x}} \int_y (h_D(\mathbf{x}) - y)^2 P(\mathbf{x}, y) P(D) d\mathbf{x} dy dD$$

$$\bar{h}(x) \triangleq E_D[h_D] = \int_D h_D p(D) dD$$

Regression with Gaussian noise

- This expression tells the expected quality of our model with random training data and a random test data

$$\begin{aligned}
 E_{\substack{(\mathbf{x}, y) \sim P \\ D \sim P^n}} [(h_D(\mathbf{x}) - y)^2] &= \int_D \int_{\mathbf{x}} \int_y (h_D(\mathbf{x}) - y)^2 P(\mathbf{x}, y) P(D) d\mathbf{x} dy dD \\
 &= E_{\mathbf{x}, y, D} [\underbrace{(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))}_{\textcircled{1}} + \underbrace{\bar{h}(\mathbf{x}) - y}_{\textcircled{2}}]^2 \\
 &= E_{\mathbf{x}, y, D} [\underbrace{(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))^2}_{\textcircled{1}} + 2 \underbrace{(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)}_{\textcircled{2}} + \underbrace{(\bar{h}(\mathbf{x}) - y)^2}_{\textcircled{3}}]
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{1} \quad E_{x|D}[(h_D(x) - h(x))^2] &= E_{x|D}[(h_D(x) - h(x))^2] \\
 &\leftarrow E_y[f(x)] \\
 &= \int f(x)p(y)dy \\
 &= f(x) \int p(y)dy \\
 &= f(x)
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \quad E_{x|D}[2(h_D(x) - h(x))(h(x) - y)] &= E_{x|D}[(h(x) - y) \underbrace{E_D[h_D(x) - h(x)]}_{h_D(x)}] \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{3} \quad E_{x|D}[(h(x) - y)^2] &= E_{x|D}[(\underbrace{h(x) - \bar{y}(x)}_{\textcircled{4}} + \underbrace{\bar{y}(x) - y}_{\textcircled{5}})^2] \\
 &= E_{x|D}[(h(x) - \bar{y}(x))^2 + 2(h(x) - \bar{y}(x))(\bar{y}(x) - y) + (\bar{y}(x) - y)^2] \\
 &= E_{x|D}[\underbrace{(h(x) - \bar{y}(x))^2}_{\textcircled{4}} + \underbrace{2(h(x) - \bar{y}(x))(\bar{y}(x) - y)}_{\textcircled{5}} + \underbrace{(\bar{y}(x) - y)^2}_{\textcircled{6}}]
 \end{aligned}$$

$\bar{y}(x) = E_{y|x}[y] = \int y p(y|x) dy$

$$⑤ E_{x,y} [2(h(x) - \bar{y}(x))(\bar{y}(x) - y)]$$

$$= \iint 2(h(x) - \bar{y}(x))(\bar{y}(x) - y) p(x, y) dx dy$$

$$= \int 2(h(x) - \bar{y}(x)) \left[\int (\bar{y}(x) - y) \overset{p(y|x)p(x)}{p(y|x)} dy \right] p(x) dx$$

$$E_{y|x} [\bar{y}(x) - y]$$

→ 0

Regression with Gaussian noise

- This expression tells the expected quality of our model with random training data and a random test data

$$E_{\substack{(\mathbf{x}, y) \sim P \\ D \sim P^n}} \left[(h_D(\mathbf{x}) - y)^2 \right] = \int_D \int_{\mathbf{x}} \int_y (h_D(\mathbf{x}) - y)^2 P(\mathbf{x}, y) P(D) d\mathbf{x} dy dD$$

$$\underbrace{E_{\mathbf{x}, y, D} \left[(h_D(\mathbf{x}) - y)^2 \right]}_{\text{Expected Test Error}} = \underbrace{E_{\mathbf{x}, D} \left[(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))^2 \right]}_{\text{Variance}} + \underbrace{E_{\mathbf{x}, y} \left[(\bar{y}(\mathbf{x}) - y)^2 \right]}_{\text{Noise}} + \underbrace{E_{\mathbf{x}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right]}_{\text{Bias}^2}$$

$$\bar{h}(\mathbf{x}) \triangleq E_D[h_D] = \int_D h_D p(D) dD$$

$$\bar{y}(\mathbf{x}) \triangleq E_{y|\mathbf{x}}[y] = \int_y y p(y|\mathbf{x}) dy$$

Expected model

Expected groundtruth

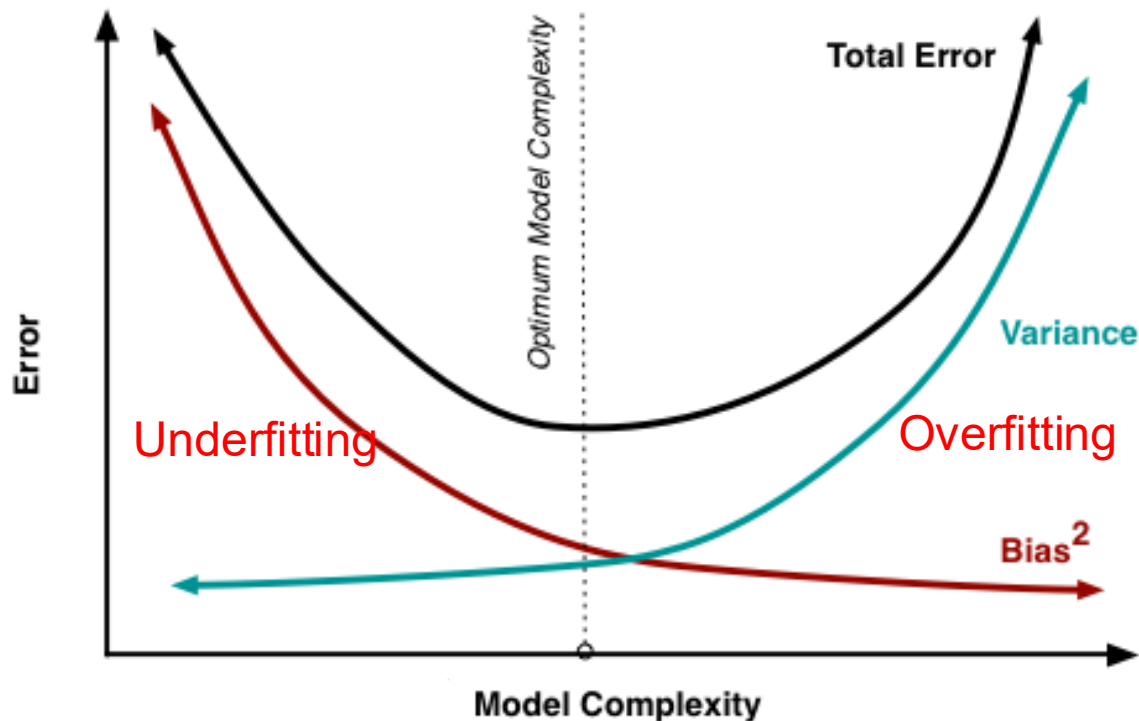
Variance, Bias, and noise

$$\underbrace{E_{\mathbf{x},y,D} \left[(h_D(\mathbf{x}) - y)^2 \right]}_{\text{Expected Test Error}} = \underbrace{E_{\mathbf{x},D} \left[(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))^2 \right]}_{\text{Variance}} + \underbrace{E_{\mathbf{x},y} \left[(\bar{y}(\mathbf{x}) - y)^2 \right]}_{\text{Noise}} + \underbrace{E_{\mathbf{x}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right]}_{\text{Bias}^2}$$

- Variance: how your classifier changes if the training data changes. Measures **generalizability**.
- Bias: The model's inherent error. If you have **infinite training** data, you will have the average classifier \bar{h} and still left with this error.
 - For example, even with infinite training data, a linear classifier will still have errors if the distribution is non-linear.
- Noise: data-intrinsic noise. Noise from measurement, noise from feature extraction, etc. Regardless of your model this remains.

Bias-Variance Underfitting-Overfitting

- *Usually* if you try to reduce the bias of your model, the variance will increase, and vice versa.
- Called the “bias-variance trade-off”

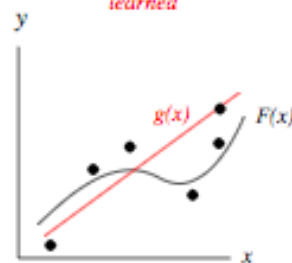
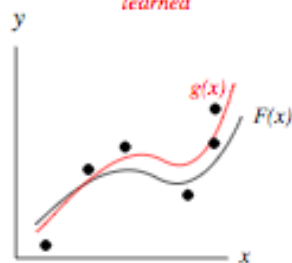
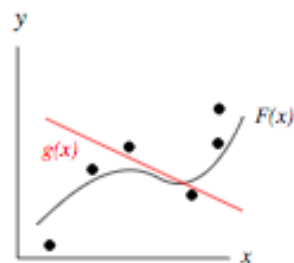
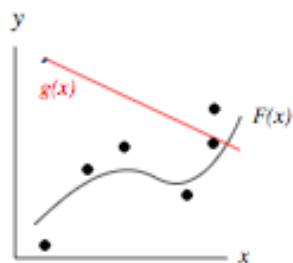
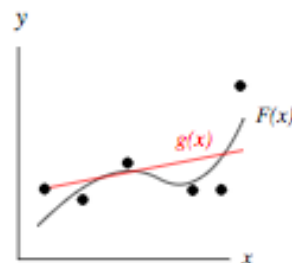
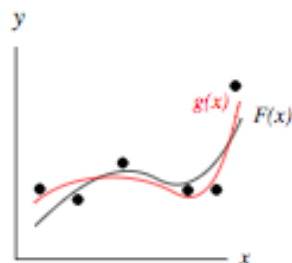
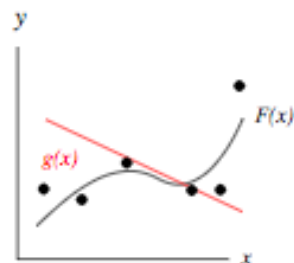
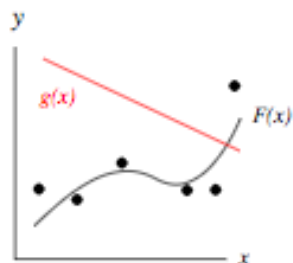
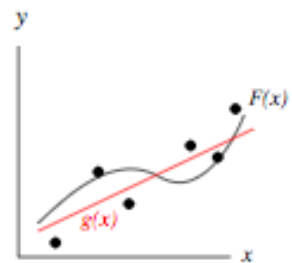
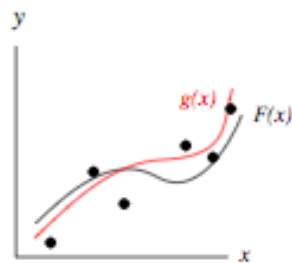
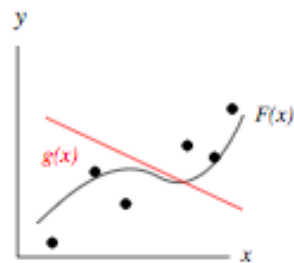
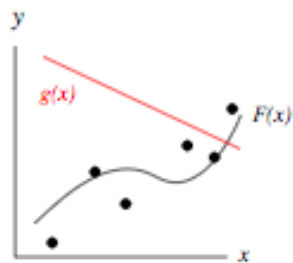


a)

b)

c)

d)

 $g(x) = \text{fixed}$ $g(x) = \text{fixed}$ $g(x) = a_0 + a_1x + a_2x^2 + a_3x^3$
learned $g(x) = a_0 + a_1x$
learned D_1  D_2  D_3 

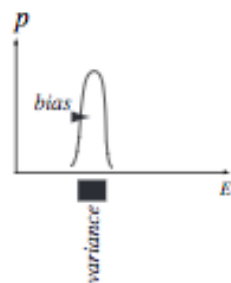
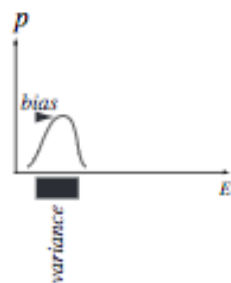
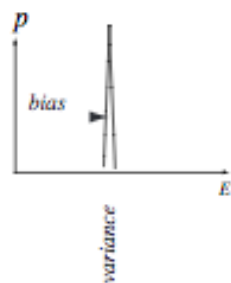
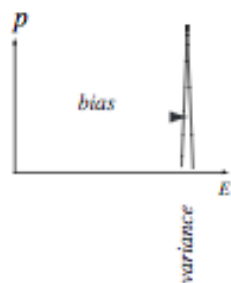
⋮

⋮

⋮

⋮

⋮



Summary

Linear Regression

Minimizing a loss function

Solve directly vs Iterative (gradient descent)

MSE view vs probabilistic view

Linear regression vs Logistic regression

For linear regression, direct or iterative should yield similar answers as long as the learning rate is small enough
(No local minima to get stuck - Convex problem)

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

$$\theta_j \Leftarrow \theta_j + r \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i) x_i^{(j)}$$

$$\theta_j \Leftarrow \theta_j + r \sum_{i=1}^m (y_i - h_{\theta}(x_i)) x_i^{(j)}$$