

# Likelihood and Naïve Bayes

---

# Predicting amount of rainfall



<https://esan108.com/%E0%B8%9E%E0%B8%A3%E0%B8%B0%E0%B9%82%E0%B8%84%E0%B8%81%E0%B8%B4%E0%B8%99%E0%B8%AD%E0%B8%B0%E0%B9%84%E0%B8%A3-%E0%B8%AB%E0%B8%A1%E0%B8%B2%E0%B8%A2%E0%B8%96%E0%B8%B6%E0%B8%87.html>

# (Linear) Regression

- $h_{\theta}(x) = \theta_0 + \theta_1 \underline{x_1} + \theta_2 \underline{x_2} + \theta_3 \underline{x_3} + \theta_4 \underline{x_4} + \theta_5 \underline{x_5}$

- $\theta$ s are the parameter (or weights)

Assume  $x_0$  is always 1

- We can rewrite

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i = \theta^T \mathbf{x}$$

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



# LMS regression with gradient descent

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


Interpretation?

# Logistic Regression

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

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

# Logistic Regression update rule

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

Update rule for linear regression

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

# Overview

L2 MSE

L1

- Probabilistic view of linear regression
  - Solution for logistic regression
- Homework 1 notes
  - Overfitting vs Underfitting (Bias – variance trade-off)
- Bayes decision models
  - Parameter estimation
    - MLE, MAP
  - Naïve Bayes

# Distribution parameter estimation

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

HHTTH

$$P(X; \theta) = P(HHTTH; \theta)$$

$$= P(H; \theta) P(H; \theta) P(T; \theta) P(T; \theta) P(H; \theta)$$

$$= (\theta)(\theta)(1-\theta)(1-\theta)(\theta)$$

$$\frac{\partial}{\partial \theta} = \theta^3(1-\theta)^2 \rightarrow \theta^3 2(1-\theta)(-1) + (1-\theta)^3 3\theta^2 = 0$$

$$\bullet L(\theta) = P(X; \theta) = P(HHTTH; \theta) \cancel{\theta^2(1-\theta)}[-2\theta + (1-\theta)3]$$

- Maximum Likelihood Estimate (MLE)

- Likelihood - Probability of encountering the data X given the parameters  $\theta$

$$2\theta = 3 - 3\theta \Rightarrow \theta = \frac{3}{5}$$

# Linear Regression Revisit

$$\cdot h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \theta_5 x_5$$

- $\theta$ s are the parameter (or weights)

- We can rewrite

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i = \theta^T \mathbf{x}$$

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



# Probabilistic Interpretation of linear regression

$$\text{Normal}(5, 1) \xrightarrow{-5} \text{Normal}(0, 1)$$

- 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  $\sigma^2$ 
  - Why Gaussian?
  - Why saying mean is zero is a valid assumption?

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

Diagram illustrating the probabilistic interpretation:

- The observed value  $y_i$  is the sum of the predicted value  $\theta^T \mathbf{x}_i$  and the error term  $\epsilon_i$ .
- The error term  $\epsilon_i$  is represented as a red circle labeled "noise".
- The predicted value  $\theta^T \mathbf{x}_i$  is represented as a red circle labeled "RV" (Random Variable).
- A red bracket labeled "noise" encompasses the error term  $\epsilon_i$ .
- A red bracket labeled "RV" encompasses the predicted value  $\theta^T \mathbf{x}_i$ .

# Probabilistic view of Linear regression

(x-11)  
26<sup>z</sup>

$$p(H|T; \theta)$$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{x^2}{2\sigma^2})$$

- Find  $\theta$
- Maximize Likelihood of seeing x and y in training

- From our assumption we know that

$$y_i = \theta^T \mathbf{x}_i + \epsilon_i \sim N(0, \sigma^2) \rightarrow N(\theta^T \mathbf{x}_i, \sigma^2)$$

error

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

Gaussian pdf

Error term is normally distributed with mean 0 and variance  $\sigma^2$

$P(y_1, y_2, y_3, y_4)$  iid

# Maximizing Likelihood

What is the assumption here?  
Is it accurate?

$\arg \max L(\theta) = \prod_{i=1}^m p(y_i | \mathbf{x}_i; \theta)$

We use the log likelihood instead  $\log(L(\theta)) = I(\theta)$

$$\log L(\theta) = \sum_{i=1}^m \log p(y_i | \mathbf{x}_i; \theta)$$

$$= \sum_{i=1}^m \left[ \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{1}{2\sigma^2} (y_i - \theta^T \mathbf{x}_i)^2 \right]$$

$$= \sum_i \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \underbrace{\sum_i \frac{1}{2\sigma^2} (y_i - \theta^T \mathbf{x}_i)^2}_{\text{error}}$$

From our previous lecture

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

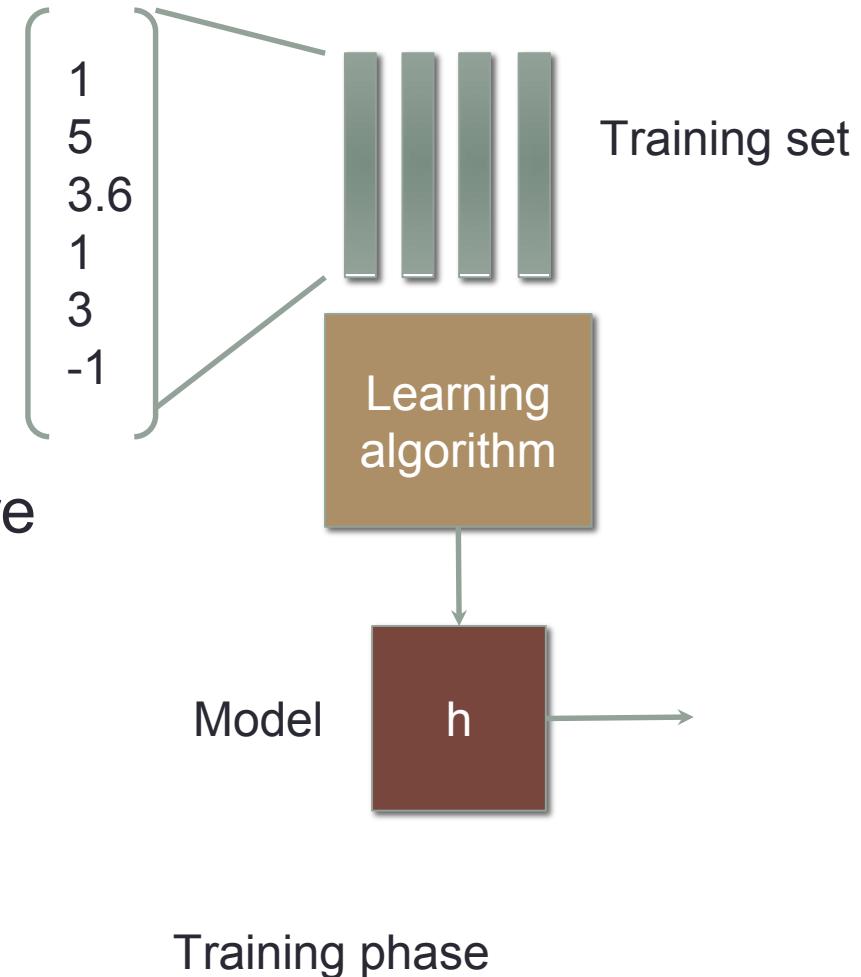
# Mean square error solution and MLE solution

- 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  $\sigma^2$  falls out of the maximization, so this derivation is true regardless of which assumption for variance is.
- 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  $\hat{y}$  is normally distributed.

# Flood or no flood

- What would be the output?
- $y = 0$  if not flooded
- $y = 1$  if flooded

0.5 1

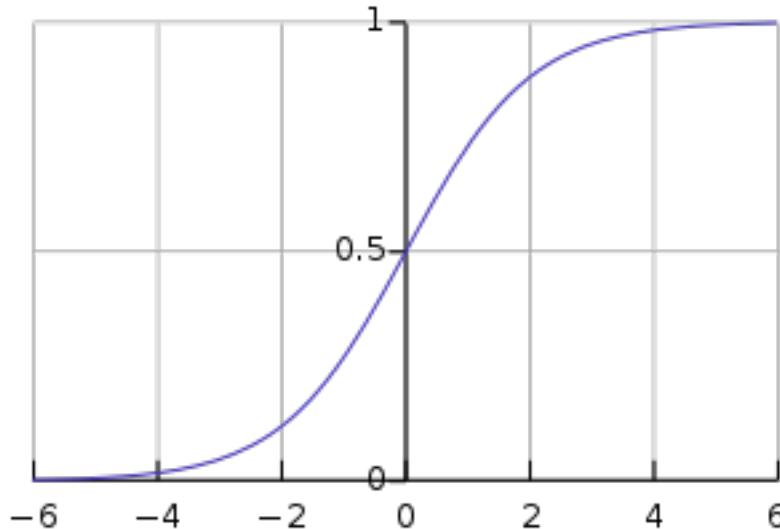


# Can we use regression?

- Yes
- $h_{\theta}(x) = \theta_0 + \theta_1x_1 + \theta_2x_2 + \theta_3x_3 + \theta_4x_4 + \theta_5x_5$
- 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)



$$f(x) = \frac{1}{1 + e^{-x}} \\ = \frac{e^x}{1 + e^x}$$

# 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

*g* ~~is~~ not logistic

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



# Probabilistic view of Logistic Regression

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

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

$$\begin{aligned} P(y=1 | x; \theta) &= h_{\theta}(x) \\ P(y=0 | x; \theta) &= 1 - h_{\theta}(x) \end{aligned}$$

or

$$\underbrace{p(y | x; \theta)}_{=} = \frac{(h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}}{1 - h_{\theta}(x)} = h_{\theta}(x)$$

# Maximizing log likelihood

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

$$\begin{aligned} L(\theta) &= P(y_1, y_2, \dots, y_m | X; \theta) \\ &\stackrel{\text{iid}}{=} \prod_{i=1}^m P(y_i | x_i; \theta) \end{aligned}$$

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

$$h_\theta(x) = \underline{g(\theta^T x)}$$

$$l(\theta) = \sum_{i=1}^m \left[ y_i \log h_\theta(x_i) + (1-y_i) \log(1-h_\theta(x_i)) \right]$$

$$\frac{\partial l(\theta)}{\partial \theta} = \sum_{i=1}^m y_i \frac{1}{h_\theta(x_i)} h'_\theta(x_i) + \sum_{i=1}^m (1-y_i) \frac{1}{1-h_\theta(x_i)} (-h'_\theta(x_i))$$

# Maximizing log likelihood

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

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

$$\frac{\partial \ell(\theta)}{\partial \theta} = \sum_{i=1}^m y_i \frac{h'_\theta(x_i)}{h_\theta(x_i)} + \sum_{i=1}^m (1-y_i) \frac{(-h'_\theta(x_i))}{1-h_\theta(x_i)}$$

$$\frac{\partial h_\theta(x_i)}{\partial \theta_j} = \frac{\partial g(\theta^T x_i)}{\partial \theta_j} = g(\theta^T x_i) (1 - g(\theta^T x_i)) \frac{\partial \theta^T x_i}{\partial \theta_j}$$

$$= \sum_{i=1}^m \left[ \frac{y_i}{h_\theta(x_i)} - \frac{(1-y_i)}{1-h_\theta(x_i)} \right] g(\theta^T x_i) (1 - g(\theta^T x_i)) x_k^{(j)}$$

$$= \sum_{i=1}^m (y_i (1-h) - (1-y_i) h) x_k^{(j)}$$

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

# 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

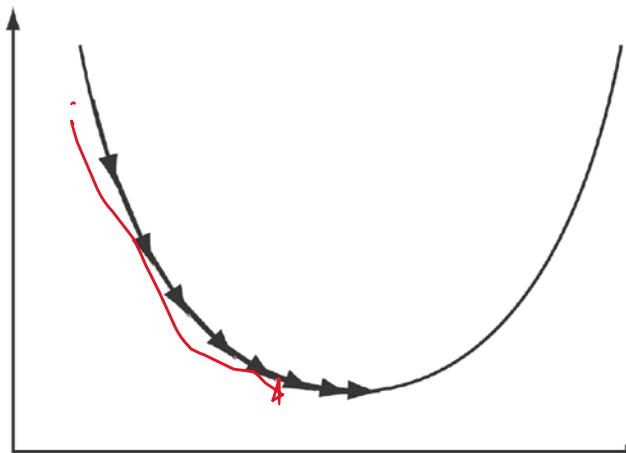
$$\theta_j \leftarrow \theta_j + r \sum_{i=1}^m (y_i - \theta^T \mathbf{x}_i) x_i^{(j)}$$

# Loose Ends from HW

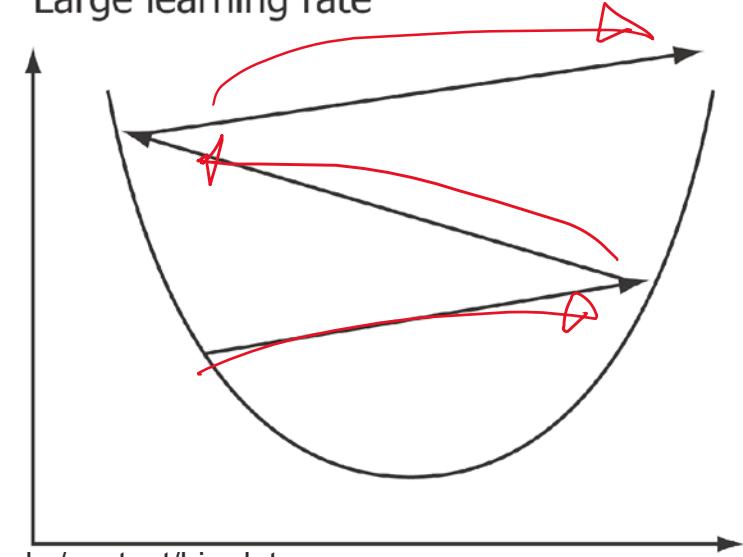
- How to select r? (The learning rate)

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

Small learning rate



Large learning rate



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

r too small and the model converges slowly

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

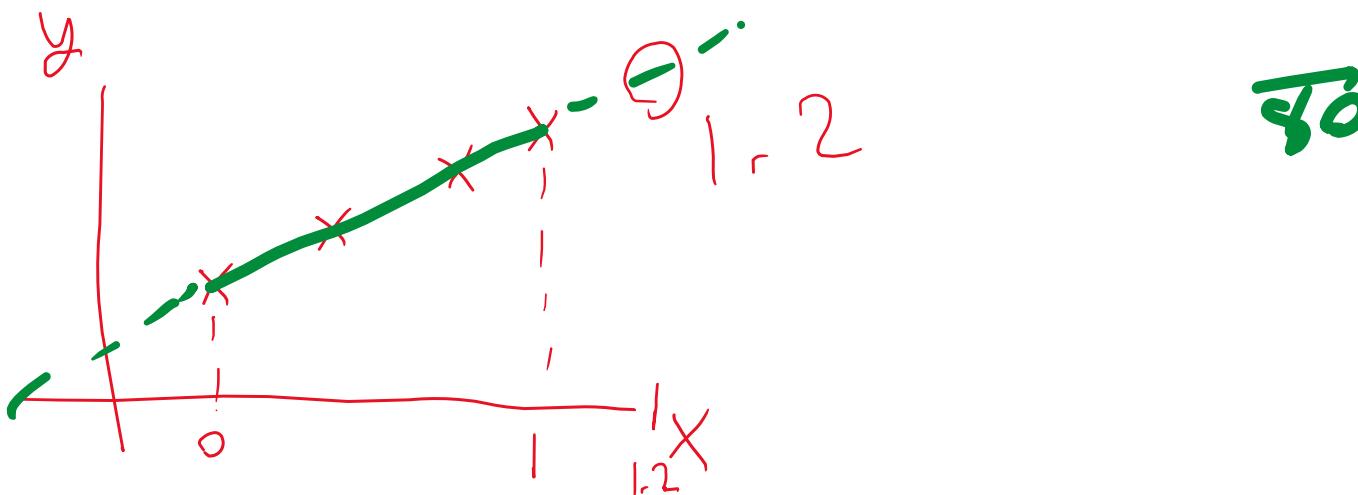
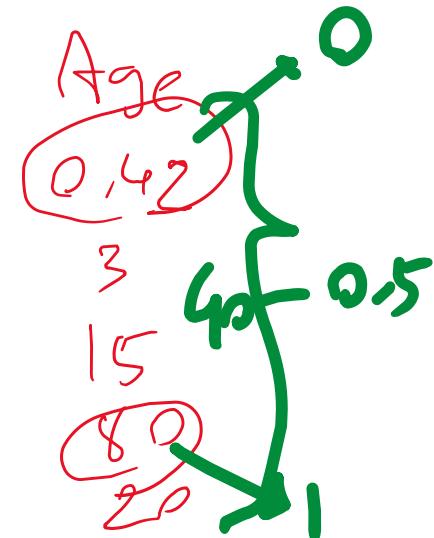
Usually have a decay over time

# 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.
- $h_{\theta}(x) = \theta_0 + \theta_1x_1 + \theta_2x_2 + \theta_3x_3 + \theta_4x_4$

# 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

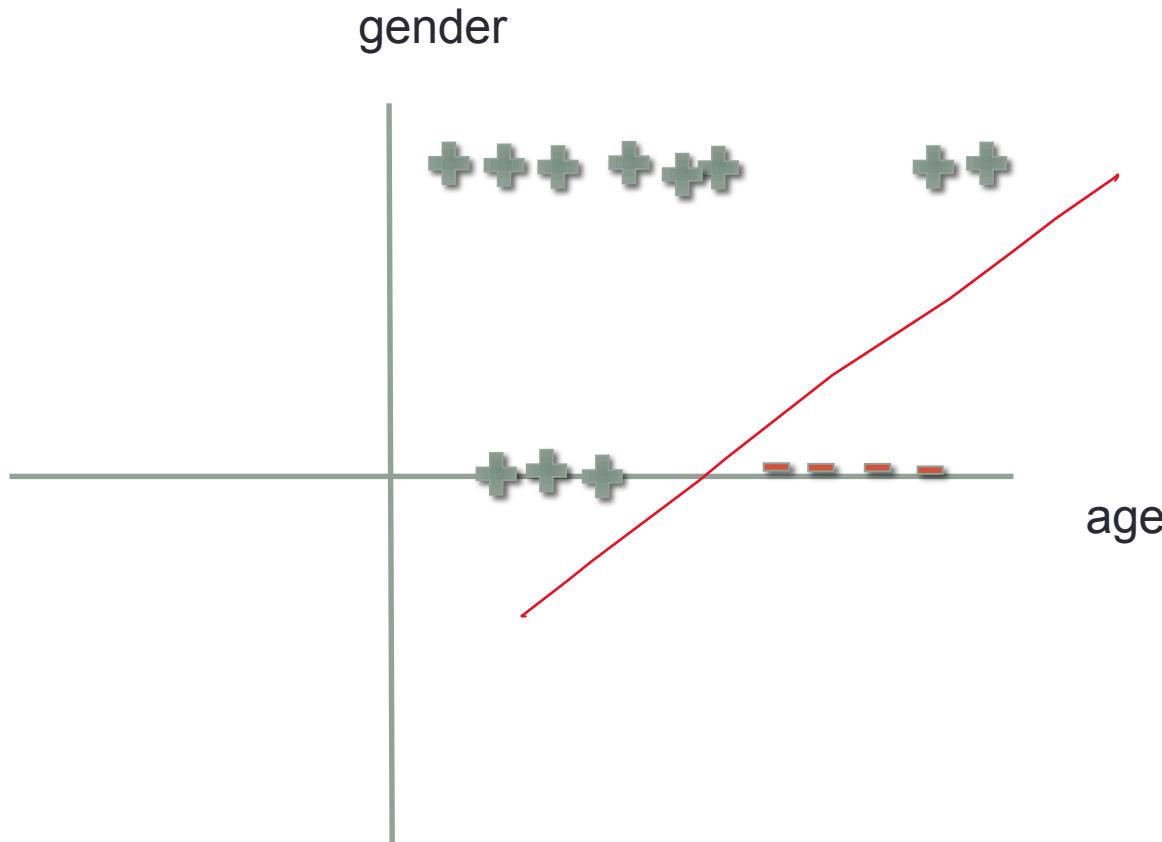


# Feature selection

- 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 ✓ + feature importance
  - Boosting ↗

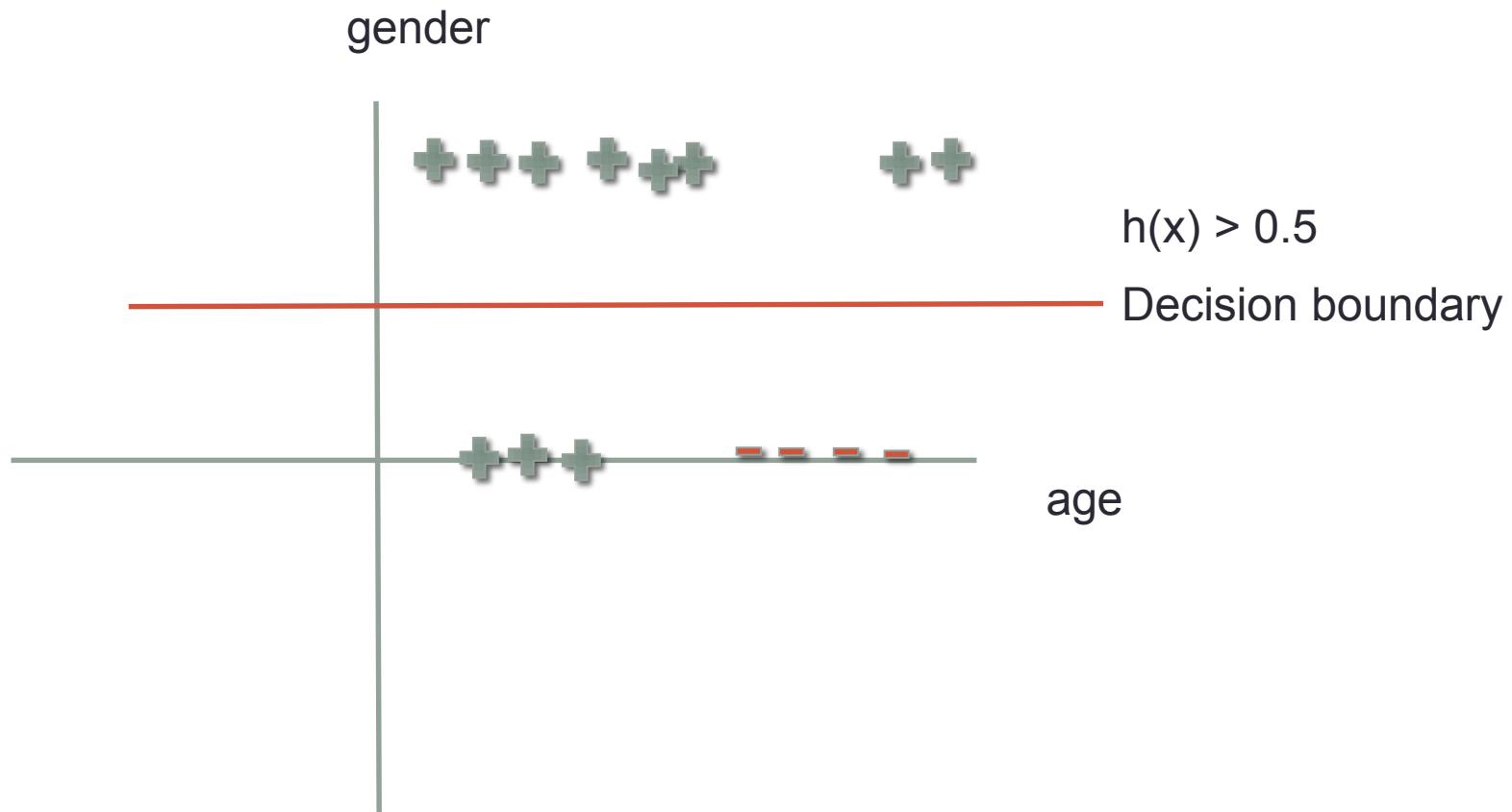
# Feature engineering

- Logistic regression is a linear classification



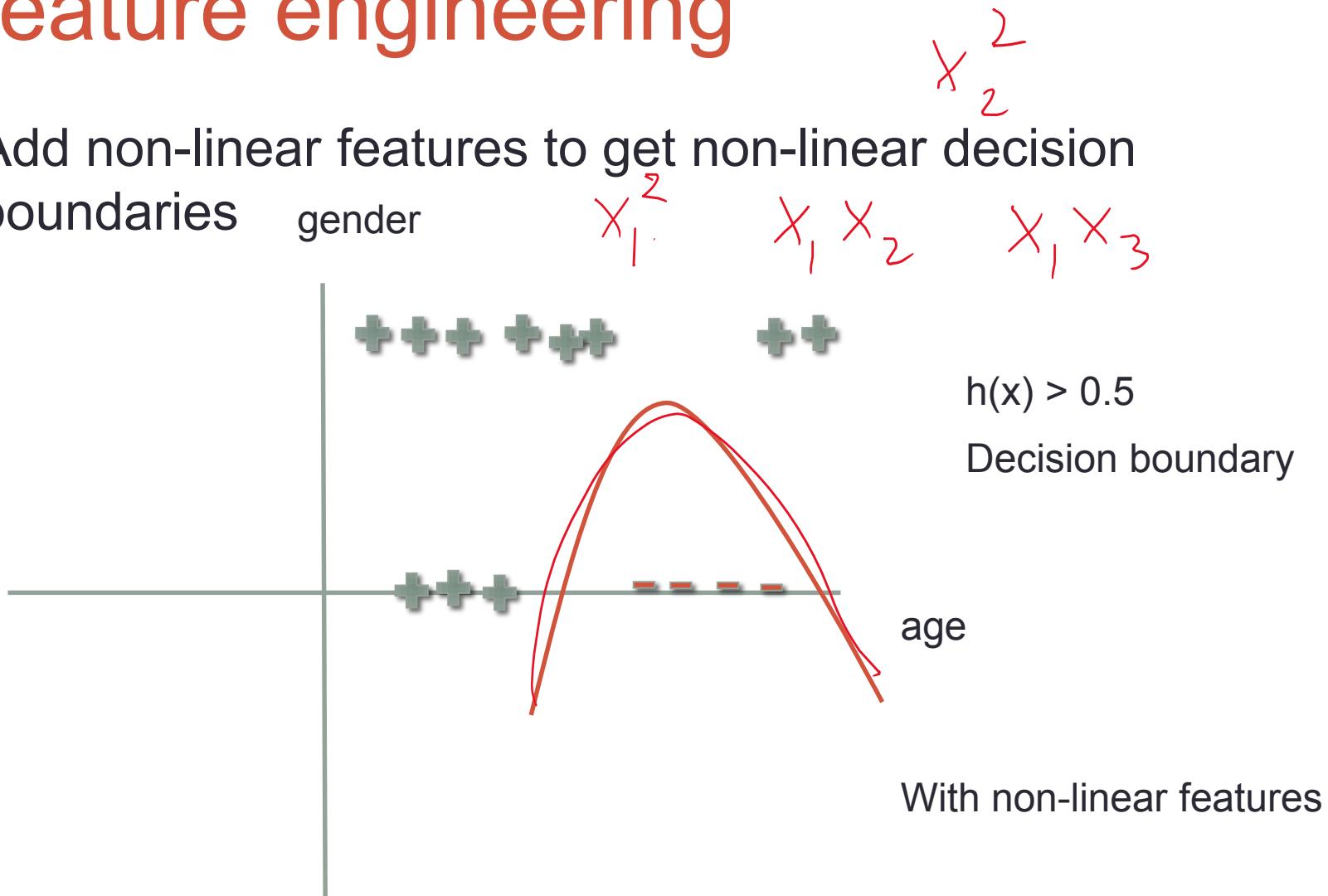
# Feature engineering

- Logistic regression is a linear classification



# Feature engineering

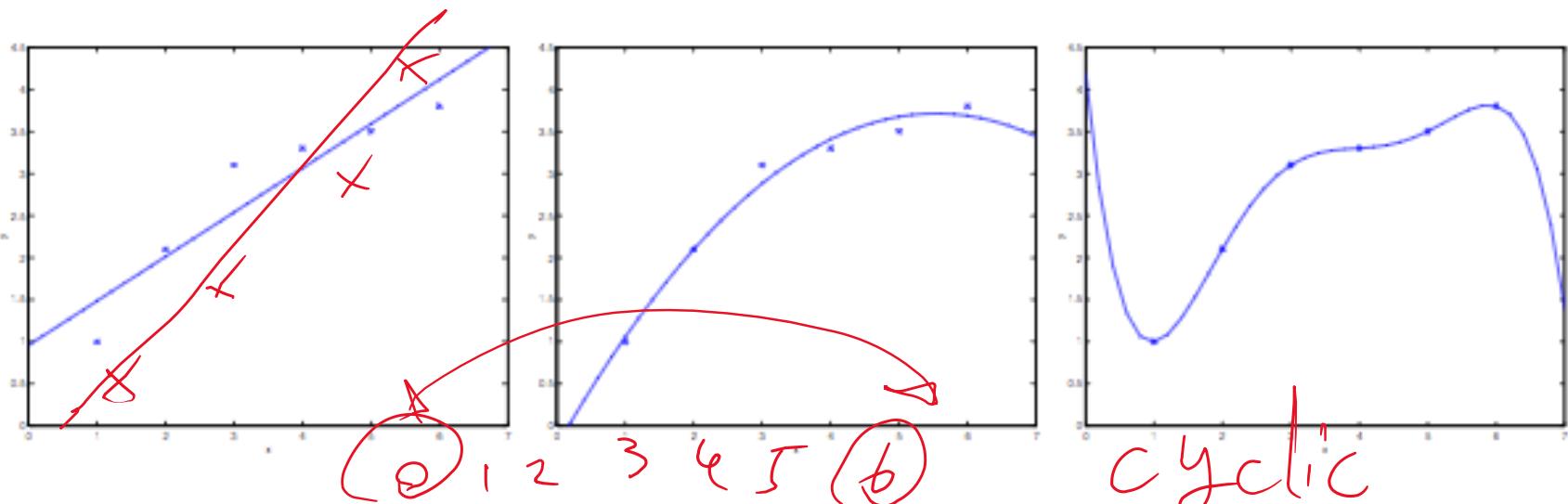
- Add non-linear features to get non-linear decision boundaries



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

14:44

# Overfitting Underfitting



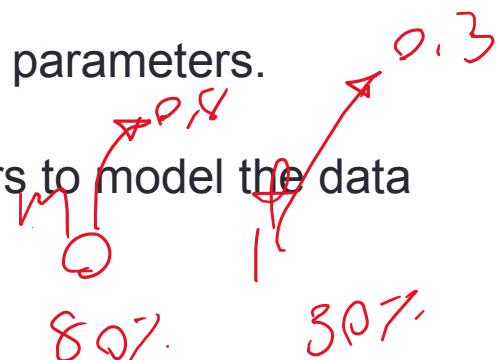
Adding more non-linear features makes the line more curvy  
(Adding more features also means more model parameters)

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**

~~target encoding~~



# Bias-Variance trade-off

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

# Regression with Gaussian noise

$$\boxed{y = h(\mathbf{x}) + \varepsilon}$$

- Where  $\varepsilon$  is normally distributed with mean zero and variance  $\sigma^2$
- The training data  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_3, y_3), (\mathbf{x}_3, y_3) \dots\}$  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  $h_D(x)$
- We calculate the expected error (squared error) on new  $(\mathbf{x}, y)$  data with the regressor

Test

- $E_{(\mathbf{x}, y)}[(h_D(\mathbf{x}) - y)^2] = \iint_{\mathbf{x} \ y} (h_D(\mathbf{x}) - y)^2 \Pr(\mathbf{x}, y) \partial y \partial \mathbf{x}$
- But  $D$  is actually random too!

$$h_D \text{ with } \Pr(\mathbf{x}, y)$$

# Regression with Gaussian noise

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

$$\cdot \mathbb{E}_{(\mathbf{x}, y)}[(h_D(\mathbf{x}) - y)^2] = \iint_{\mathbf{x} \ y} (h_D(\mathbf{x}) - y)^2 \Pr(\mathbf{x}, y) \partial y \partial \mathbf{x}$$

- 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

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

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

## Regression with Gaussian noise

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

$$\begin{aligned}
 & E_{\substack{(x,y) \sim P \\ D \sim P^n}} [(h_D(x) - y)^2] = \int_D \int_x \int_y (h_D(x) - y)^2 P(x, y) P(D) \partial x \partial y \partial D \\
 & \approx E_{X,y,D} \left[ [h_p(x) - \bar{h}(x)] + [\bar{h}(x) - y]^2 \right] \\
 & \approx E_{X,y,D} \left[ (h_p(x) - \bar{h}(x))^2 + 2(h_p(x) - \bar{h}(x))(\bar{h}(x) - y) + (\bar{h}(x) - y)^2 \right]
 \end{aligned}$$

(1) (2) (3) (4)

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

$$\textcircled{1} \quad E_{xyD} \left[ (h_p(x) - h(x))^2 \right]$$

$$= E_{xD} \left[ (h_p(x) - h(x))^2 \right]$$

$$E_y[h(x)]$$

$$= \int f(x)p(y)dy$$

$$= f(x) \int p(y)dy$$

$$= f(x)$$

$$\textcircled{2} \quad E_{xyD} \left[ 2(h_p(x) - h(x))(h(x) - y) \right]$$

$$= E_{xy} \left[ (h(x) - y) E_D[h_p(x) - h(x)] \right]$$

$$= \textcircled{3}$$

$$\textcircled{3} \quad E_{xyD} \left[ (h(x) - y)^2 \right]$$

$$= E_{xyD} \left[ (h(x) - \bar{y}(x) + \bar{y}(x) - y)^2 \right]$$

$$= E_{xyD} \left[ (h(x) - \bar{y}(x))^2 + 2(h(x) - \bar{y}(x))(\bar{y}(x) - y) + (\bar{y}(x) - y)^2 \right]$$

$$= \textcircled{4}$$

$$+ \textcircled{5}$$

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

$$= \int y p(y|x) dy$$

$$= \textcircled{6}$$

$$\textcircled{5} \quad E_{x,y|P} \left[ 2(t(x) - \bar{y}(x))(\bar{y}(x) - y) \right]$$

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


 $p(y|x)p(x)$

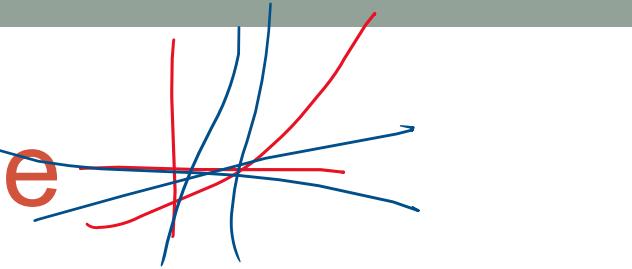
$$= \int 2(t(x) - \bar{y}(x)) \left[ \int (\bar{y}(x) - y) p(y|x) dy \right] p(x) dx$$


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

 $\bar{y}$

- Q

# Variance, Bias, and noise

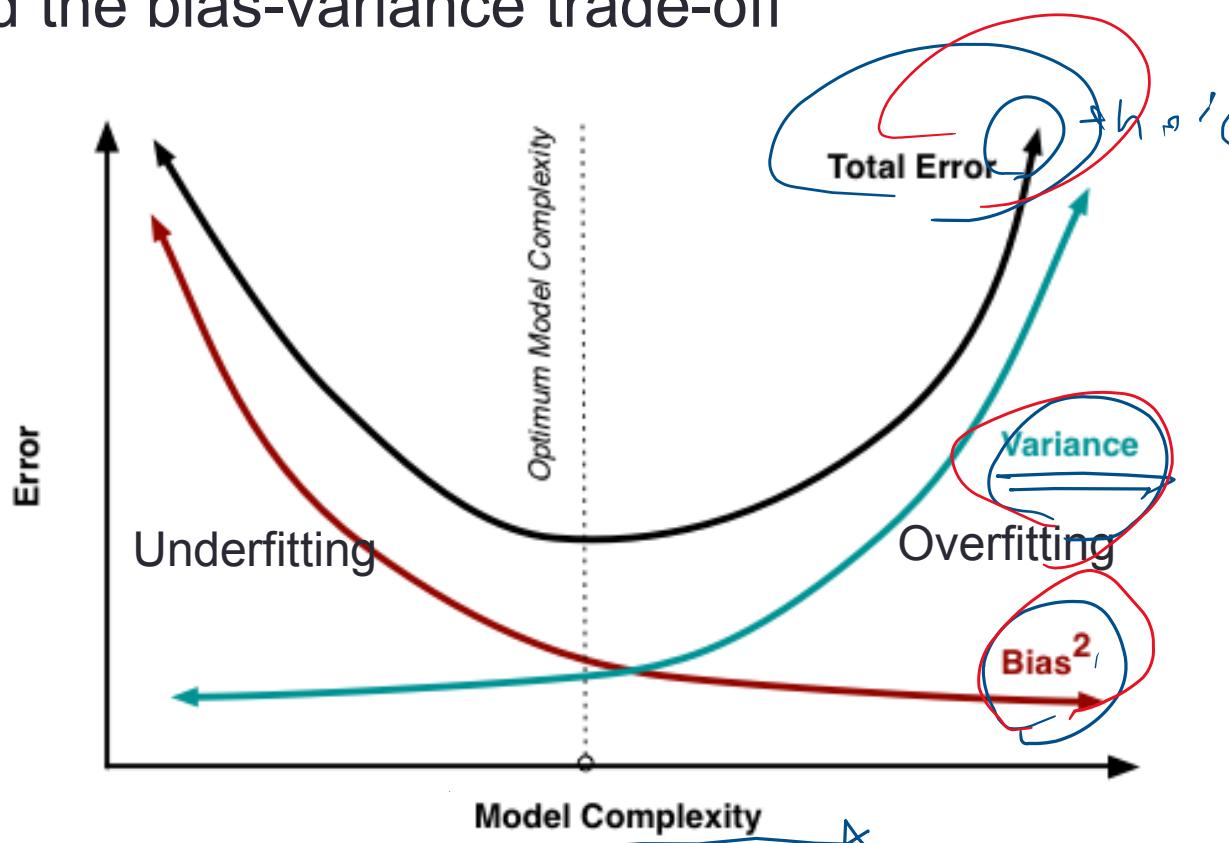


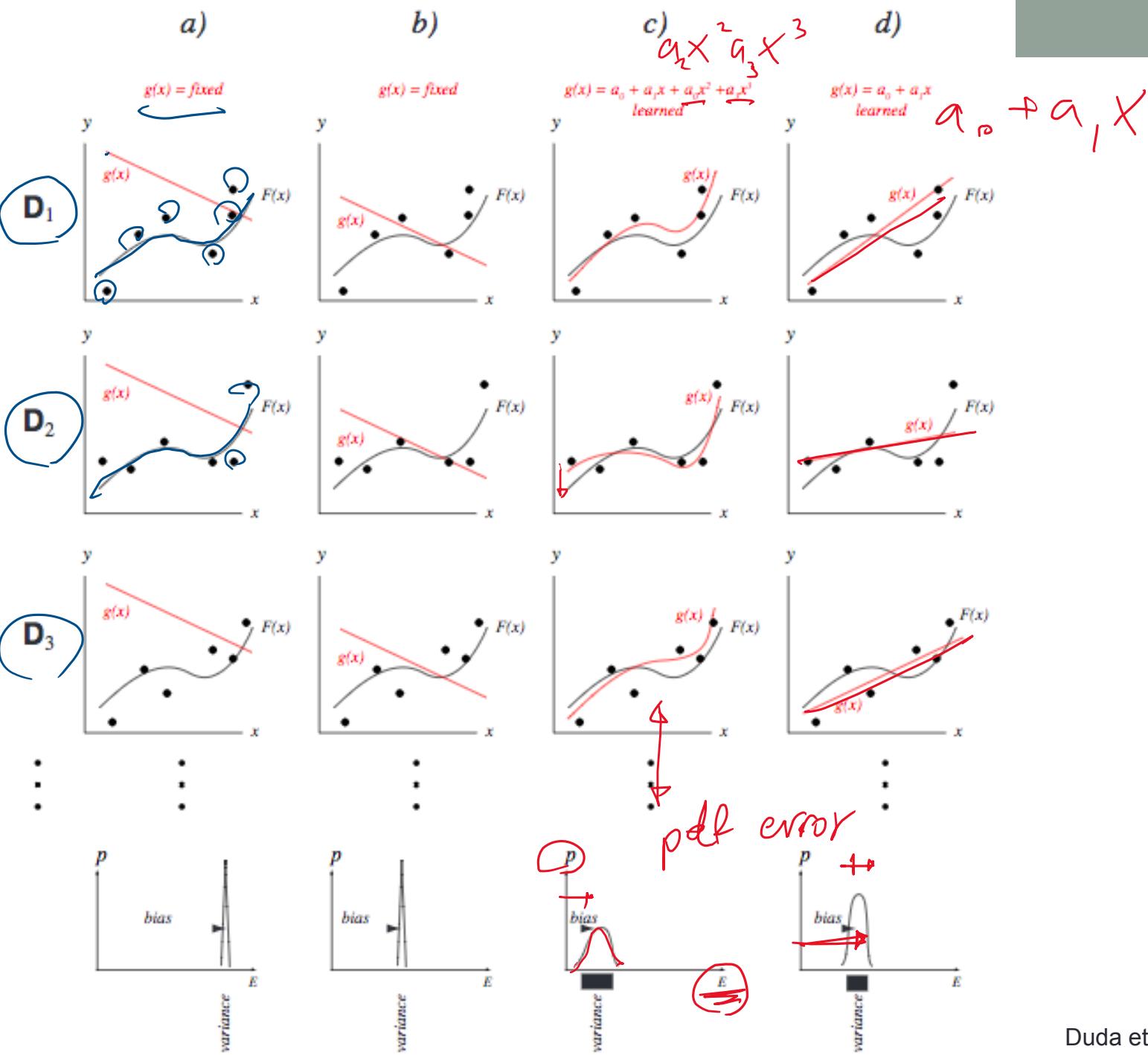
$$E_{\mathbf{x},y,D} \left[ (h_D(\mathbf{x}) - y)^2 \right] = \underbrace{E_{\mathbf{x},D} \left[ (h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))^2 \right]}_{\text{Variance}} + \underbrace{E_{\mathbf{x},y} \left[ (\bar{y}(\mathbf{x}) - y)^2 \right]}_{\text{Noise}} + \underbrace{E_{\mathbf{x}} \left[ (\bar{h}(\mathbf{x})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

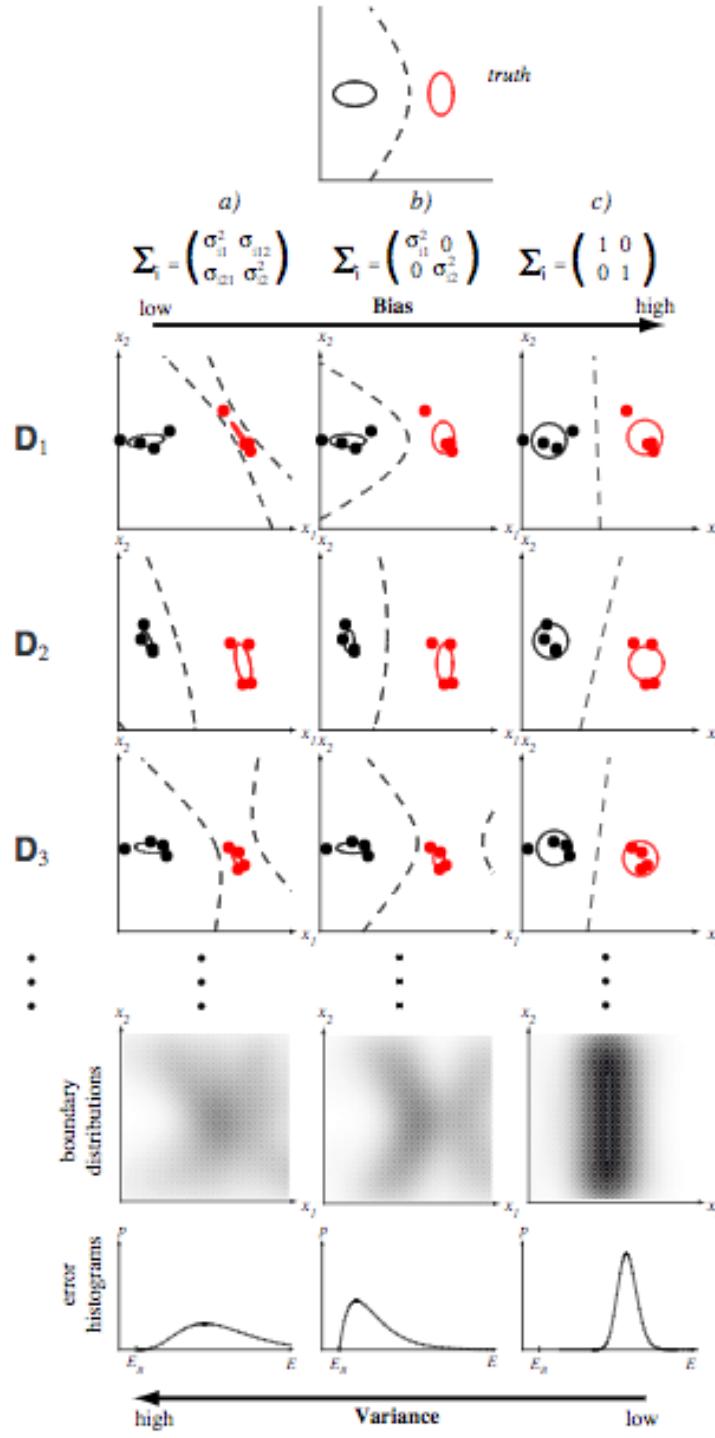
- Variance: how your classifier changes if the training data changes. Measures **generalizability**.
- Bias: The model's inherent error. If you have **infinite training** data, you will have the average classifier  $\bar{h}$  and still left with this error.
  - For example, even with infinite training 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



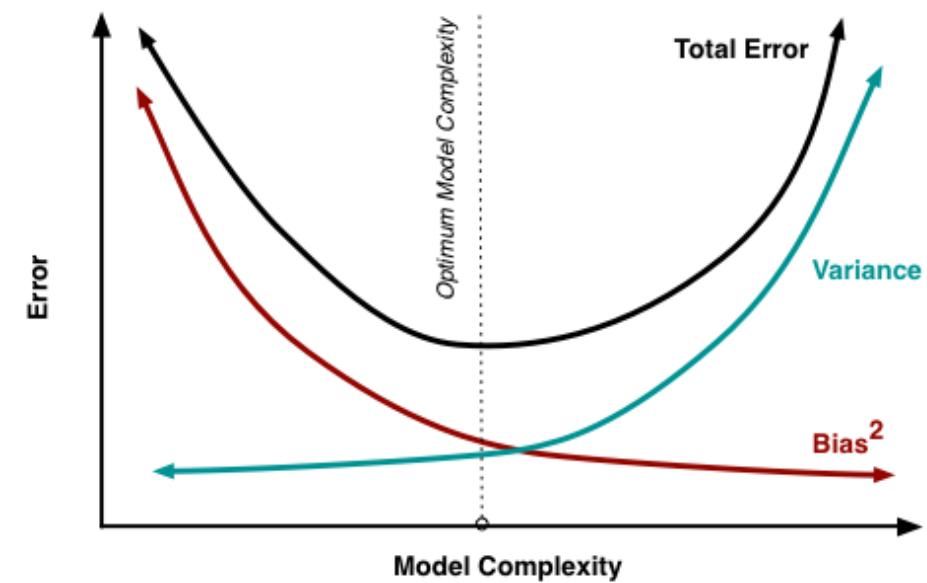




# When to stop the update?

- 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
<b>Your Best Entry ↑</b>						
Your submission scored 0.81340, which is not an improvement of your best score. Keep trying!						
453	▼ 7	Clement Sengelen		0.80861	11	2mo

# The Bayes Lecture

- Bayes Decision Rule
- Naïve Bayes

# A simple decision rule

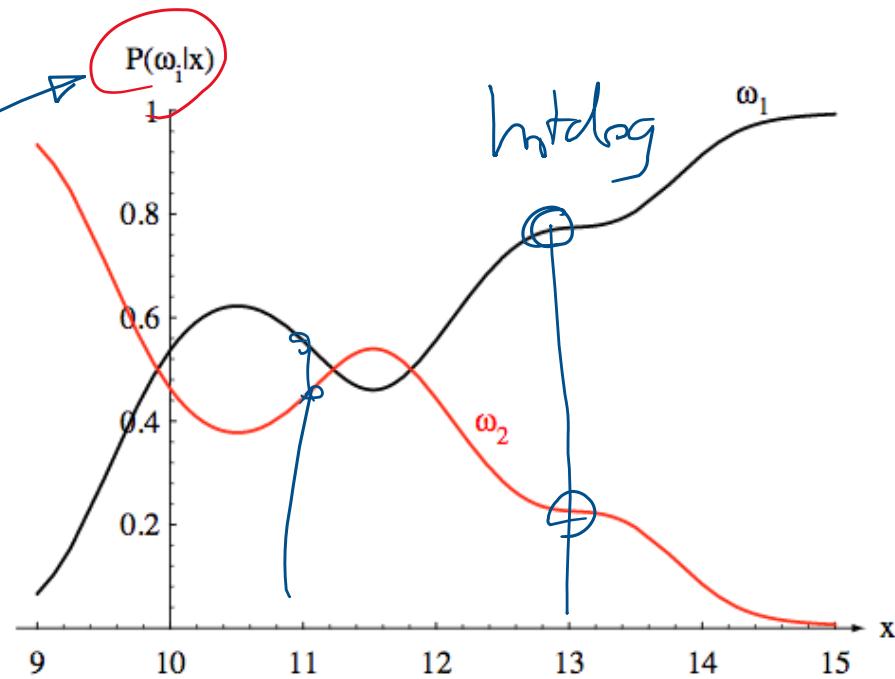
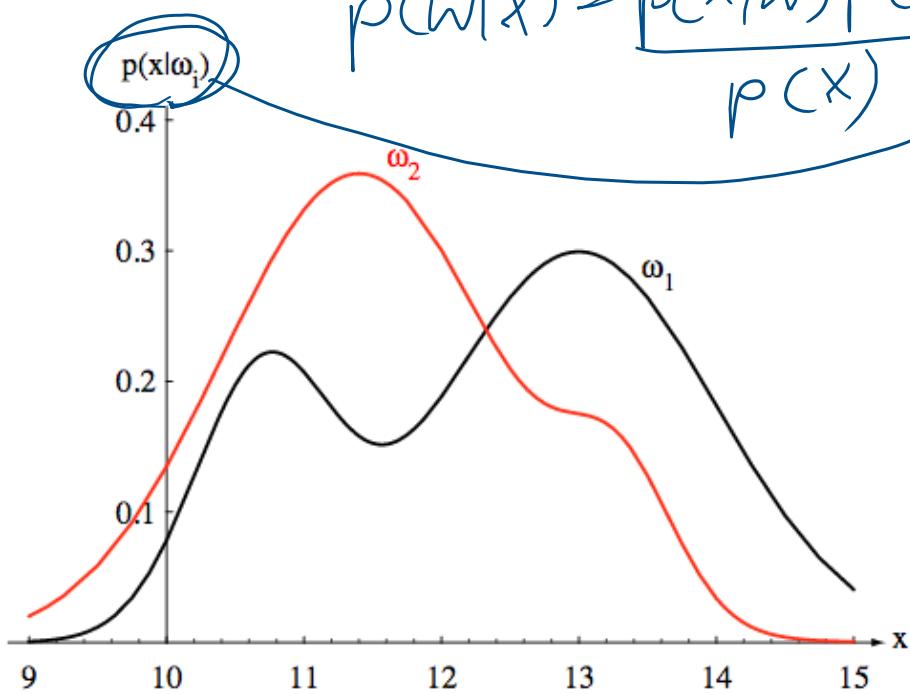
X

$\omega_1$   
hotdog  
 $\omega_2$   
not hotdog

- If we can know either  $p(x|w)$  or  $p(w|x)$  we can make a classification guess

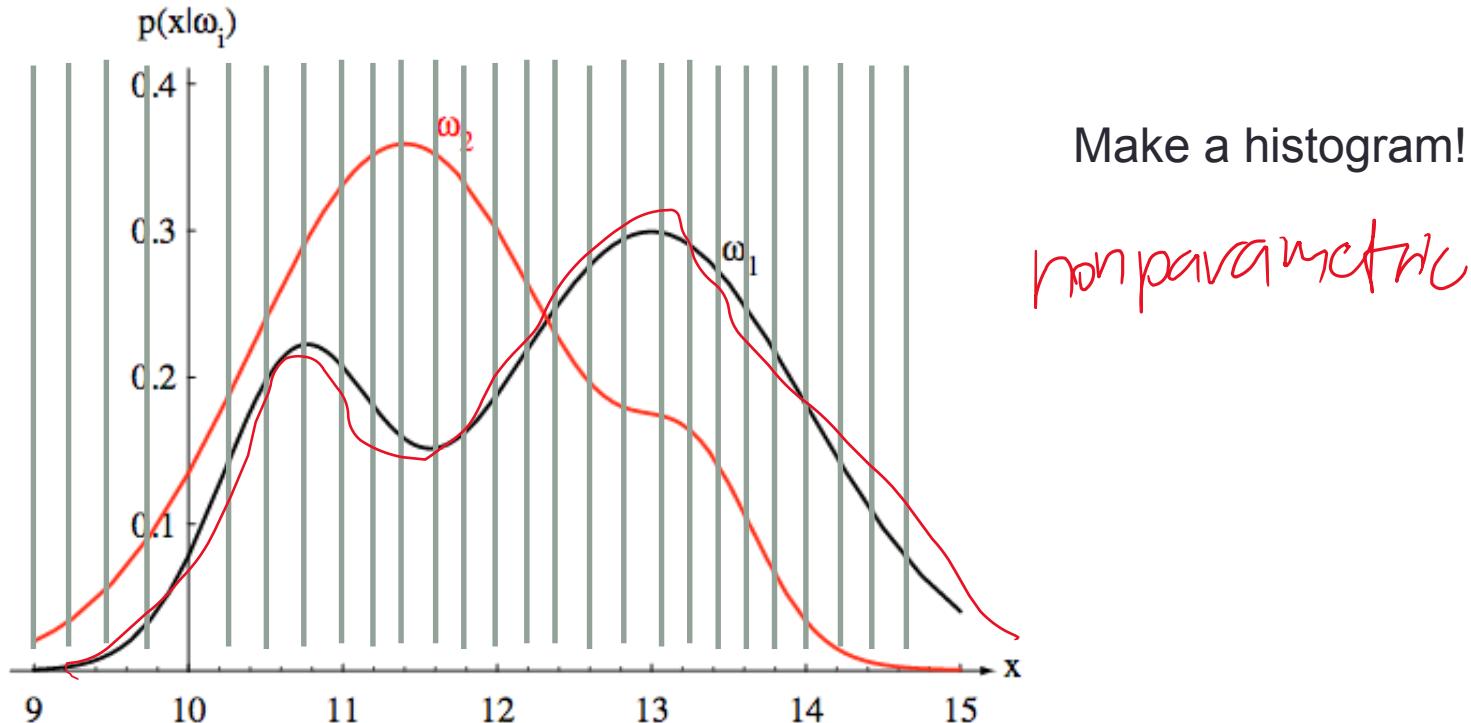
$$p(\omega_i|x) = \frac{p(x|\omega_i)p(\omega_i)}{p(x)}$$

$x = 13$



Goal: Find  $p(x|w)$  or  $p(w|x)$

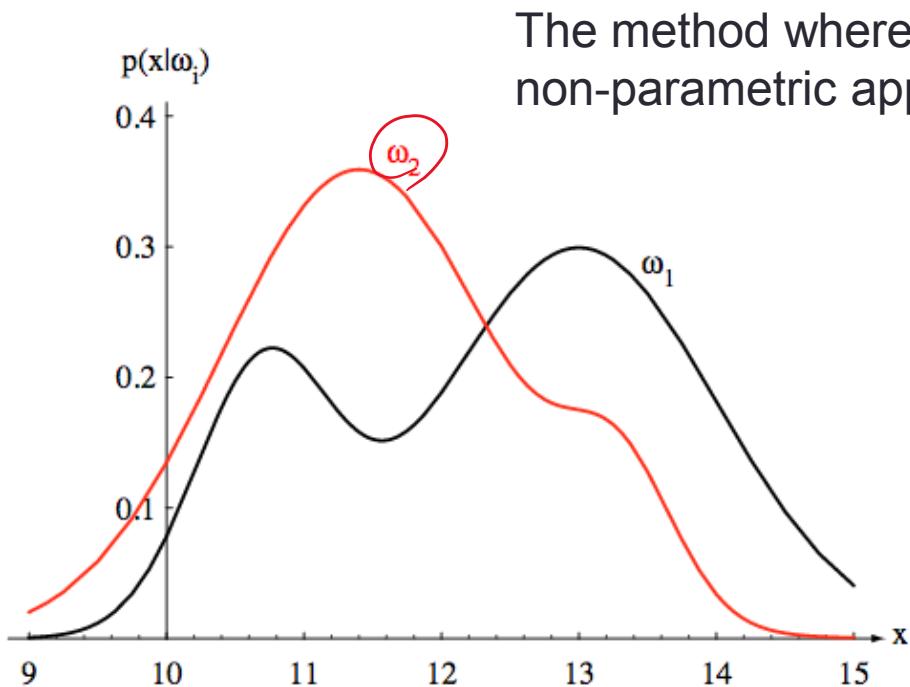
# A simple way to estimate $p(x|w)$



What happens if there is no data in a bin?

# The parametric approach

- We **assume**  $p(x|w)$  or  $p(w|x)$  follow some distributions with parameter  $\theta$



The method where we find the histogram is a non-parametric approach

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

Goal: Find  $\theta$  so that we can estimate  $p(x|w)$  or  $p(w|x)$

# Maximum Likelihood Estimate (MLE)

$$p(w_i|x) = \frac{p(x|w_i)p(w_i)}{p(x)}$$

- Maximizing the likelihood (probability of data given model parameters)

Posterior = likelihood \* prior  
evidence

$p(x|\theta) = L(\theta)$  <- This assumes the data is fixed

- Usually done on log likelihood

- Take the partial derivative wrt to  $\theta$  and solve for the  $\theta$  that maximizes the likelihood

# MLE of binomial trials

- A coin with bias is tossed N times. k times are heads. Find  $\theta$ , the probability of the coin landing head.

$$L(\theta) = \binom{N}{k} \theta^k (1-\theta)^{N-k}$$

log

$$\frac{d}{d\theta}$$

$$\theta = \frac{k}{N}$$

# MLE of Gaussian

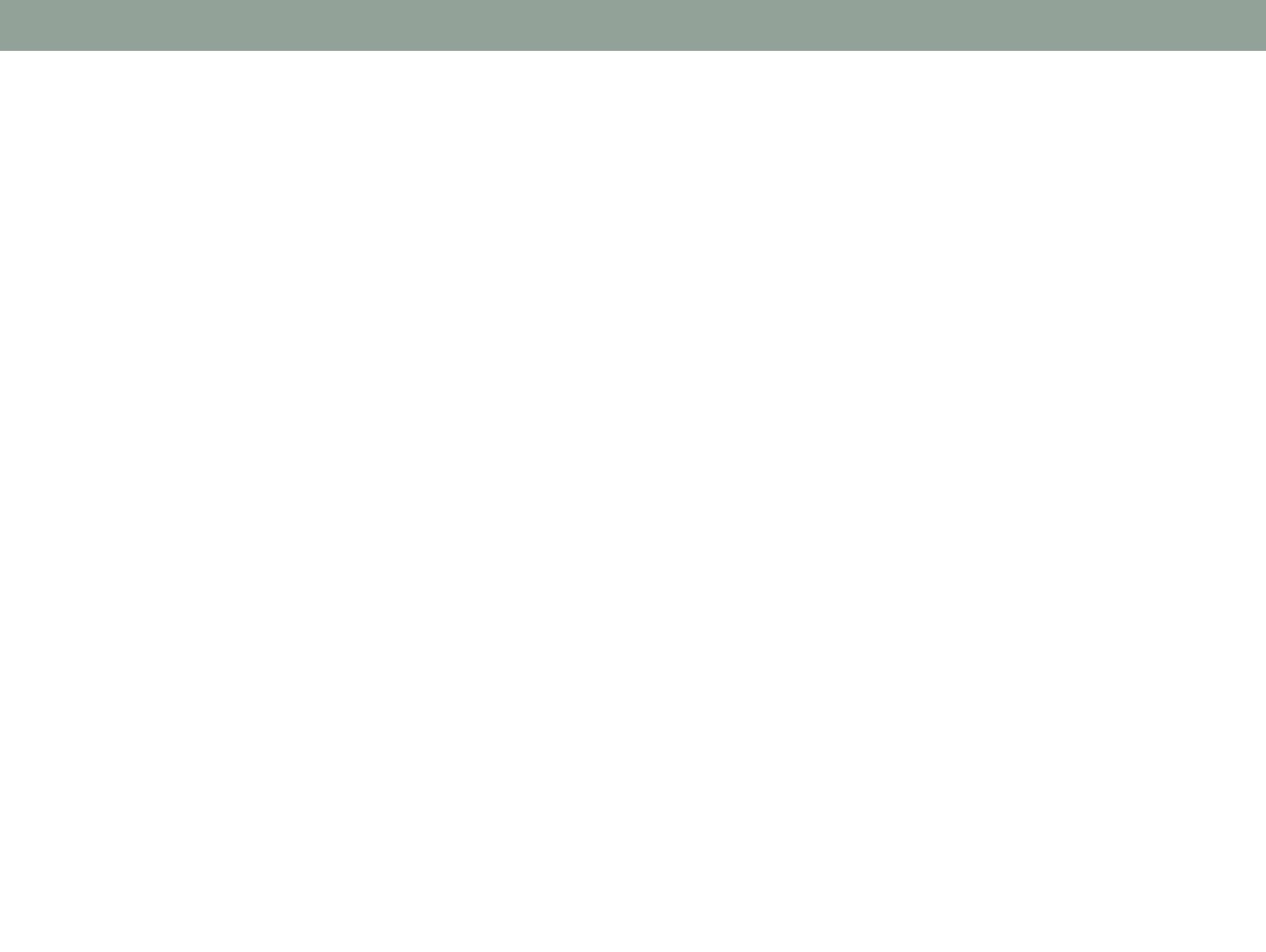
$$\mu = \theta_1 \quad \sigma^2 = \theta_2$$

- Observe  $x_i$ , estimate the mean and the variance. Assume the data is normally distributed.

$$x_1, \dots, x_n$$

$$P(x_1, \dots, x_n; \theta_1, \theta_2) = \prod_{i=1}^n P(x_i; \theta_1, \theta_2)$$

$$\frac{\partial}{\partial \theta_1} \log \sum \log \frac{1}{\sqrt{2\pi\theta_2}} + \log \exp\left(-\frac{1}{2} \frac{(x_i - \theta_1)^2}{\theta_2}\right)$$
$$\theta_1 = \frac{1}{n} \sum_{i=1}^n x_i$$
$$\theta_2 = \frac{\sum (x_i - \theta_1)^2}{n}$$



# Maximum Likelihood Estimate (MLE)

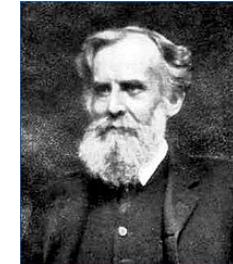
$$p(w_i|x) = \frac{p(x|w_i)p(w_i)}{p(x)}$$

Posterior = likelihood \* prior  
evidence

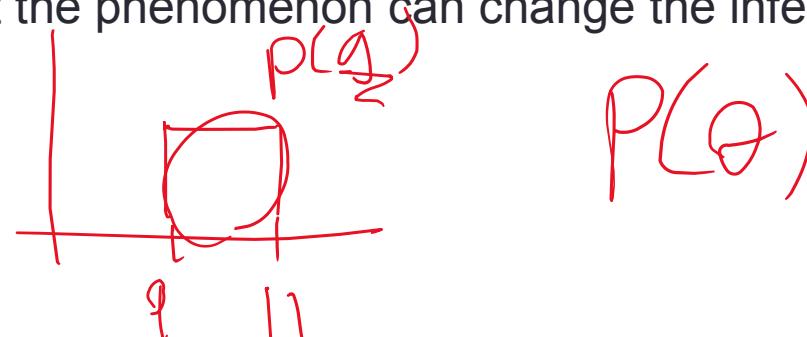
- Maximizing the likelihood (probability of data given model parameters)  
 $p(\mathbf{x}|\theta) = L(\theta)$  <- This assumes the data is fixed
- Usually done on log likelihood
- Take the partial derivative wrt to  $\theta$  and solve for the  $\theta$  that maximizes the likelihood

# Frequentist vs Bayesian view

- Frequentist  $\underset{\text{MLE}}{P(\text{data} ; \theta)}$ 
  - Probability is “frequency of occurrence”
  - Data is from a random procedure that draw from unknown but fixed phenomenon.
    - Distribution parameter is a constant
- Bayesian
  - Probability is “degree of uncertainty”
  - Data is fixed and you want to infer about the unknown phenomenon.  $\theta$ 
    - Distribution parameter is a distribution
    - Prior knowledge about the phenomenon can change the inference results.



$$g = 9,8$$



# Maximum A Posteriori (MAP) Estimate

*Freq*  
MLE

*Bayesian*  
MAP

- Maximizing the likelihood (probability of data given model parameters)

$$\underset{\theta}{\operatorname{argmax}} p(\mathbf{x}|\theta)$$

$$p(\mathbf{x}|\theta) = L(\theta)$$

- Usually done on log likelihood

- Take the partial derivative wrt to  $\theta$  and solve for the  $\theta$  that maximizes the likelihood

- Maximizing the posterior (model parameters given data)

$$\underset{\theta}{\operatorname{argmax}} p(\theta|\mathbf{x})$$

- But we don't know  $p(\theta|\mathbf{x})$

- Use Bayes rule

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$

- Taking the argmax for  $\theta$  we can ignore  $p(\mathbf{x})$

$$\underset{\theta}{\operatorname{argmax}} p(\mathbf{x}|\theta) p(\theta)$$

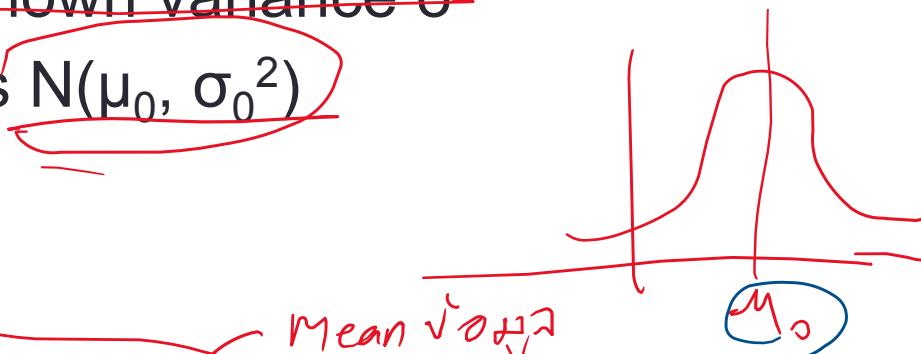
*prior*

# MAP on Gaussian

$x_1, \dots, x_n$

$$\frac{3}{5}a + \frac{2}{5}b$$

- We know  $x$  is Gaussian with unknown mean  $\mu$  that we need to estimate and known variance  $\sigma^2$
- Assume the prior of  $\mu$  is  $N(\mu_0, \sigma_0^2)$



- MAP estimate of  $\mu$  is

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \left[ \frac{1}{n} \sum_{i=1}^n x_i \right] + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

MLE estimate

11 98

argmax  $\prod_i P(X_i | \mu, \sigma^2) P(\mu)$

$$= \prod_i \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(X_i - \mu)^2}{2\sigma^2}\right) \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left(-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right)$$

log

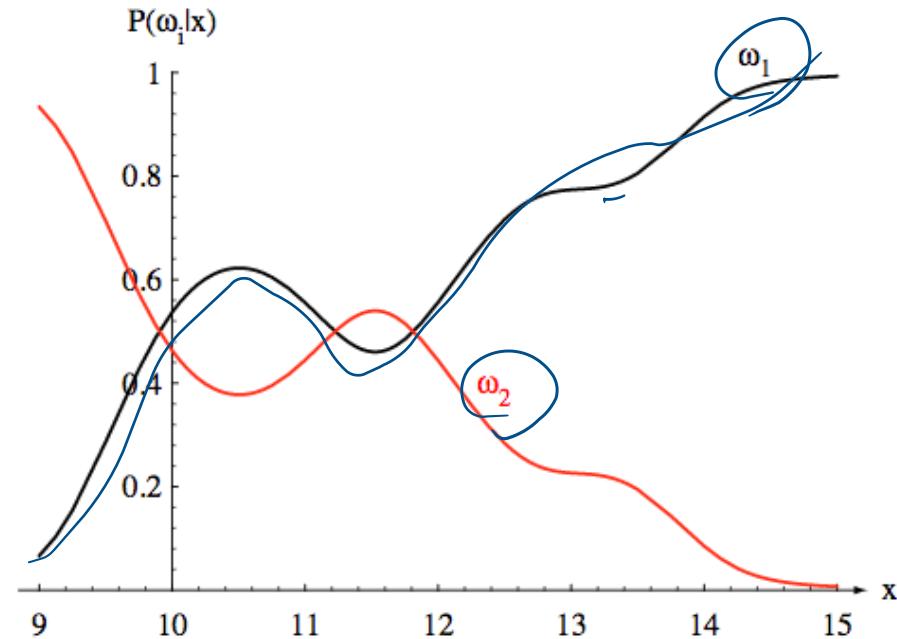
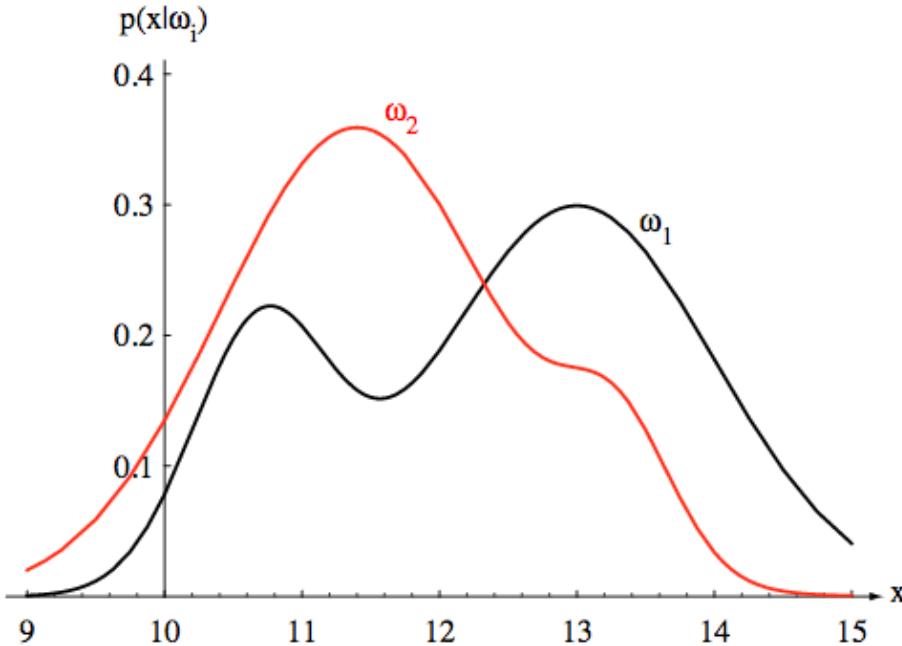
2

# Notes of MAP estimate

- Usually harder to estimate than MLE
- If we use an uninformative prior for  $\theta$ 
  - MAP estimate = MLE
- Given infinite data
  - MAP estimate converges to MLE
- MAP is useful when you have less data, so you need additional knowledge about the domain
  - MAP estimate tends to converge faster than MLE even with an arbitrary distribution
  - Can help prevent overfitting
- **Useful for model adaptation** (MAP adaptation)
  - Learn MLE on larger dataset, use this as your prior distribution
  - Learn MAP estimate on your dataset

# A simple decision rule

- If we can know either  $p(x|w)$  or  $p(w|x)$  we can make a classification guess



Goal: Find  $p(x|w)$  or  $p(w|x)$  by finding the parameter of the distribution

# Likelihood ratio test

- If  $P(w_1|x) > P(w_2|x)$ , that  $x$  is more likely to be class  $w_1$

more

- Again we know  $P(x|w_1)$  is more intuitive and easier to calculate than  $P(w_1|x)$

- Our classifier becomes

$$P(x|w_1)P(w_1) \quad ? \quad P(x|w_2)P(w_2)$$

$$\frac{P(w_1|x)}{P(x|w_1)P(w_1)} \quad ? \quad \frac{P(w_2|x)}{P(x|w_2)P(w_2)}$$

- 

$$\frac{P(x|w_1)}{P(x|w_2)}$$

?

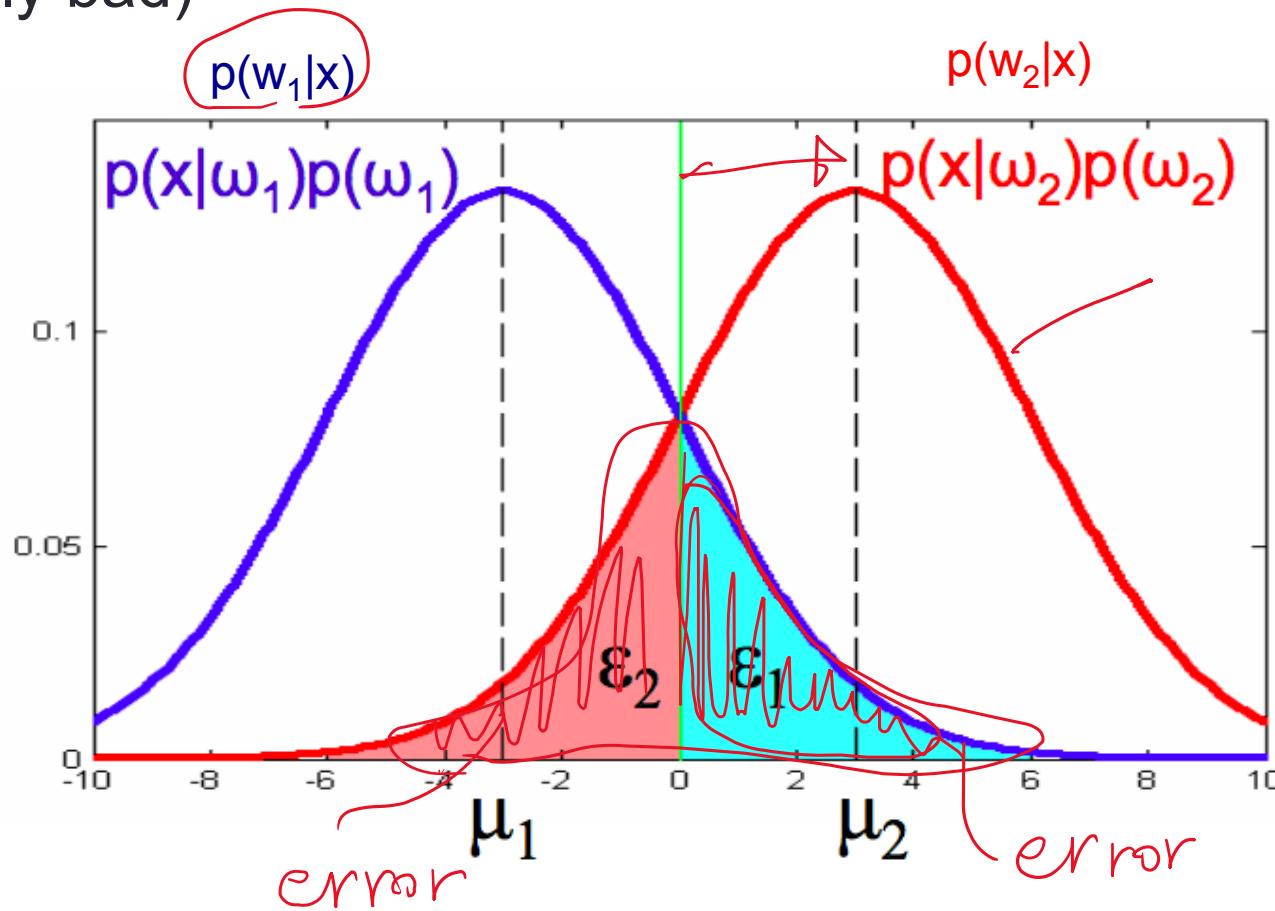
$$\frac{P(w_2)}{P(w_1)}$$

Ratio of priors

Likelihood ratio

# Notes on likelihood ratio test (LRT)

- LRT minimizes the classification error (all errors are equally bad)



# Notes on LRT

$$1 - \text{ກຳນົດ}$$
$$2 - \text{ຈຳອັບ}$$

error ຈຳອັບ

$\text{ກຳນົດ} \rightarrow \text{ຈຳອັບ}$

- If  $P(w_1|x) > P(w_2|x)$ ,  $x$  is more likely to be class  $w_1$ 
  - Also known as MAP decision rule
  - The classifier is sometimes called the Bayes classifier
- If we do not want to treat all error equally, we can assign different loss to each error, and minimize the expected loss. This is called Bayes loss/risk classifier

$$\frac{P(x|w_1)}{P(x|w_2)}$$

$$? \quad \frac{P(w_2)(L_{1|2} - L_{2|2})}{P(w_1)(L_{2|1} - L_{1|1})}$$

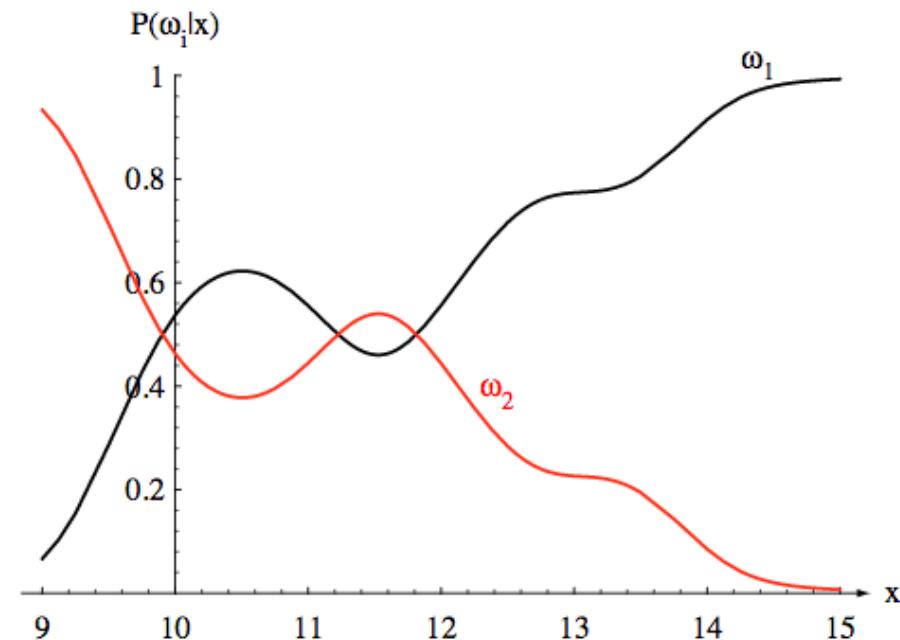
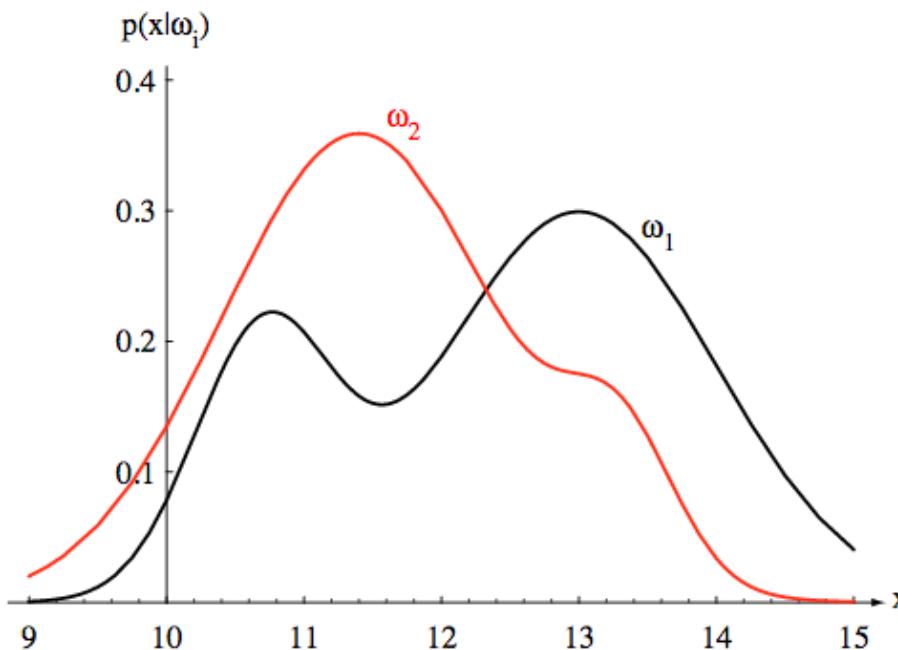
- When we treat errors equally, we refer to the zero-one loss
- $L_{1|2} = 1, L_{2|2} = 0, L_{2|1} = 1, L_{1|1} = 0$

# Notes on LRT

- If we treat the priors as equal, we get the **maximum likelihood criterion**

- $\frac{P(x|w_1)}{P(x|w_2)}$  ?

1



# Naïve Bayes

- Below is the LRT or the Bayes classifier

$$P(x|w_1)P(w_1) \quad ? \quad P(x|w_2)P(w_2)$$

- What about Naïve Bayes?

- Here  $x$  is a vector with  $m$  features  $[x_1, x_2, \dots, x_m]$
- $P(x|w_i)$  is  $m+1$  dimensional
  - Sometimes too hard to model, not enough data, overfit, *curse of dimensionality*, etc.
- Assumes  $x_1, x_2, \dots, x_m$  independent given  $w_i$  (conditional independence)
  - What does this mean?

# Modeling distributions

input

Wind in the morning

$$X \in \{\text{Calm}, \text{Windy}\}$$

output

PM2.5 level in the afternoon  $Y \in \{\text{Low}, \text{Med}, \text{High}\}$

$$\operatorname{argmax} P(Y | X) = \operatorname{argmax} P(X|Y) P(Y)$$

Day	X	Y
1	W	M
2	C	M
3	W	M
4	W	H
5	C	L
6	W	L
7	C	H
8	W	L

# Modeling distributions

Wind in the morning

$$X \in \{\text{Calm}, \text{Windy}\}$$

PM2.5 level in the afternoon

$$Y \in \{\text{Low}, \text{Med}, \text{High}\}$$

$\operatorname{argmax} P(Y | X)$

Day	X	Y
1	W	M
2	C	M
3	W	M
4	W	H
5	C	L
6	W	L
7	C	H
8	W	L

P(X, Y)	L	M	H
C	1/8	1/8	1/8
W	2/8	2/8	1/8

Joint distribution

MLE

P(Y   X)	L	M	H
C	1/3	1/3	1/3
W	2/5	2/5	1/5

Conditional distribution

9 W → L

# Modeling distributions

Wind in the morning

$X \in \{\text{Calm}, \text{Windy}\}$

PM2.5 level in the afternoon

$Y \in \{\text{Low}, \text{Med}, \text{High}\}$

$\operatorname{argmax} P(Y | X)$

Joint distribution

Day	X	Y
1	W	M
2	C	M
3	W	M
4	W	H
5	C	L
6	W	L
7	C	H
8	W	L

P(X, Y)	L	M	H
C			
W			

Total data
8

count(X, Y)	L	M	H
C	1	1	1
W	2	2	1

$$P(X, Y) = \frac{\text{Count}(X, Y)}{\text{Total count}}$$

is the Maximum Likelihood Estimate (MLE) of  $P(X, Y)$

# Modeling distributions

Wind in the morning

$$X \in \{\text{Calm}, \text{Windy}\}$$

PM2.5 level in the afternoon

$$Y \in \{\text{Low}, \text{Med}, \text{High}\}$$

$\operatorname{argmax} P(Y | X)$

Day	X	Y
1	W	M
2	C	M
3	W	M
4	W	H
5	C	L
6	W	L
7	C	H
8	W	L

P(Y   X)	L	M	H
C			
W			

Conditional distribution

Total data
8

count(X,Y)	L	M	H	Total
C	1	1	1	3
W	2	2	1	5

$P(Y | X) = \frac{\text{Count}(X, Y)}{\text{Total count}(X)}$  is the Maximum Likelihood Estimate (MLE) of  $P(Y|X)$

# Curse of dimensionality

Input

Wind in the morning

$X \in \{\text{Calm, Windy}\}$

PM2.5 level in the afternoon

$Y \in \{\text{Low, Med, High}\}$

Output  
PM2.5 level in the evening

$Z \in \{\text{Low, Med, High}\}$

$$\operatorname{argmax} P(Z | Y, X) = \operatorname{argmax} P(Y, X | Z) P(Z)$$

Day	X	Y	Z
1	W	L	M
2	C	M	M
3	W	H	M
4	W	M	H
5	C	M	L
6	W	M	L
7	C	L	H
8	W	H	L

count(Z, Y, X)	Z=L	Z=M	Z=H
X=W, Y=L	0	1	0
X=W, Y=M	1	0 or 2	1
X=W, Y=H	1	1	0
X=C, Y=L	0	0	1
X=C, Y=M	1	1	0
X=C, Y=H	0	0	0

(8)

# Naïve Bayes

$$\bullet P(x|w_i)P(w_i) = P(w_i) \prod_j P(x_j|w_i)$$

- This assumption simplifies the calculation

$$P(x_1, x_2, x_3 | w_i) = P(x_1 | w_i) P(x_2 | w_i) P(x_3 | w_i)$$

Conditiona | Independent

$$x_1 \perp x_2 | w$$

# Simplifying assumptions

Wind in the morning

$$X \in \{\text{Calm}, \text{Windy}\}$$

PM2.5 level in the afternoon

$$Y \in \{\text{Low}, \text{Med}, \text{High}\}$$

PM2.5 level in the evening

$$Z \in \{\text{Low}, \text{Med}, \text{High}\}$$

$$\begin{aligned} \operatorname{argmax} P(Z | Y, X) &= \operatorname{argmax} P(Y, X | Z) P(Z) \\ &= \operatorname{argmax} P(Y|Z)P(X|Z)P(Z) \end{aligned}$$

Day	X	Y	Z
1	W	L	M
2	C	M	M
3	W	H	M
4	W	M	H
5	C	M	L
6	W	M	L
7	C	L	H
8	W	H	L

q

$P(Y   Z)$	$Y = L$	M	H
$Z = L$			
M			
H			

$P(X   Z)$	$X = W$	C
$Z = L$		
M		
H		

# Dealing with zero probs

1. Use a very small value instead of zero (~~flooring~~)
2. Smooth the values using counts from other observations (smoothing)
3. Use priors (MAP adaptation)

Day	X	Y	Z
1	W	L	M
2	C	M	M
3	W	H	M
4	W	M	H
5	C	M	L
6	W	M	L
7	C	L	H
8	W	H	L

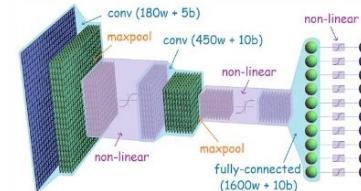
- 20  
| P

P(Y   Z)	Y = L	M	H
Z = L	0		
M			
H			

P(X   Z)	X = W	C
Z = L		
M		
H		

# WHO WOULD WIN?

AN INCREDIBLY COMPLEX  
MULTI-LAYER CONVOLUTIONAL  
NEURAL NETWORK



ONE NAIVE BOI



# Naïve Bayes Notes

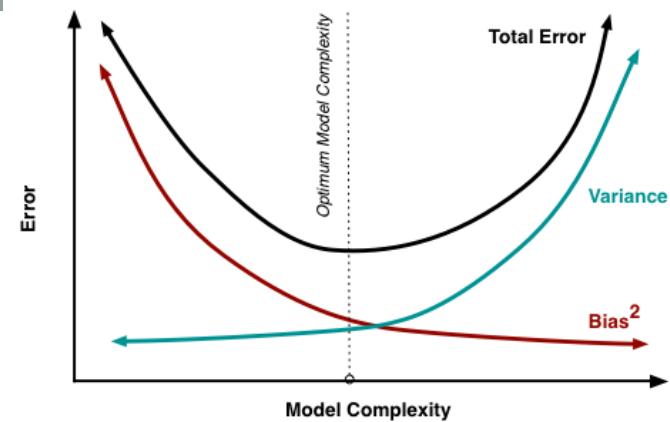
$$\bullet P(x|w_i)P(w_i) = P(w_i) \prod_j P(x_j|w_i)$$

- Note that we do not say anything about what kind of distribution  $P(x_j|w_i)$  is.  ~~$P(\text{apple}|w_i)$~~  ?  ~~$P(\text{apple}|w_0)$~~   $P(\text{apple}/w_0)$   $P(\text{apple}/w_1)$
- In the homework you will play with this
  - Clean data
  - Estimate  $P(x_j|w_i)$  using MLE, parametric and non-parametric version
  - Do prediction
  - Understand more about metrics
- Naïve Bayes can handle missing data
- Naïve Bayes is fast and quite good in practice
  - [https://www.reddit.com/r/datascience/comments/hmhg9v/why\\_is\\_naive\\_bayes\\_so\\_popular\\_for\\_nlp/](https://www.reddit.com/r/datascience/comments/hmhg9v/why_is_naive_bayes_so_popular_for_nlp/)

# Next homework

# Summary

- Probabilistic view of linear regression
- Bias-Variance trade-off
  - Overfitting and underfitting
- MLE vs MAP estimate
  - How to use the prior
- LRT (Bayes Classifier)
  - Naïve Bayes



$$\frac{P(x|w_1)}{P(x|w_2)}$$

Likelihood ratio

?

$$\frac{P(w_2)}{P(w_1)}$$

Ratio of priors

