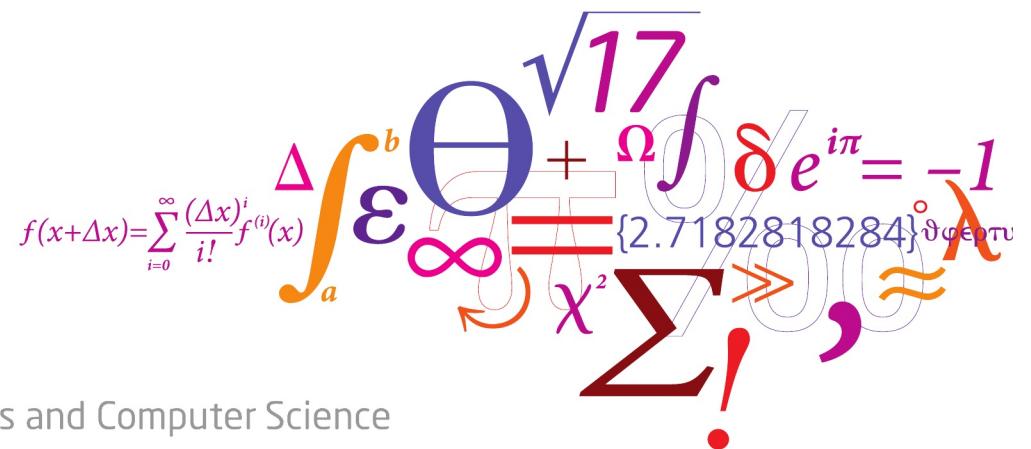


02450: Introduction to Machine Learning and Data Mining

Artificial Neural Networks and Bias/Variance

Tue Herlau

DTU Compute, Technical University of Denmark (DTU)

$$f(x+\Delta x) = \sum_{i=0}^{\infty} \frac{(\Delta x)^i}{i!} f^{(i)}(x)$$


DTU Compute

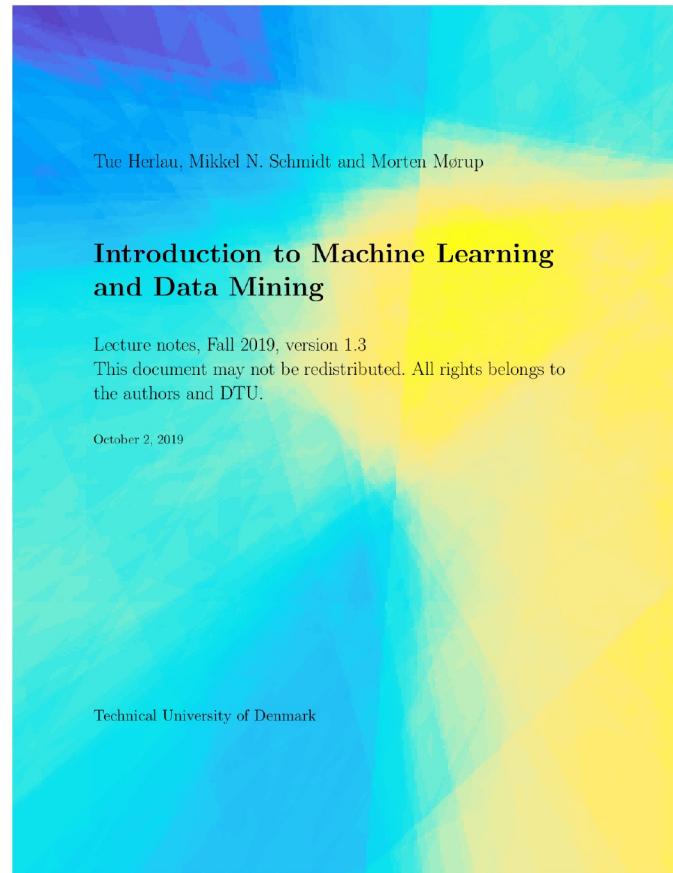
Department of Applied Mathematics and Computer Science

Today

Feedback Groups of the day:

Julie Liu Olin, Nicklas Leander Lund, Jakob Andreas Lodberg, Asbjørn Kjær Olling, David Collin Jørgensen, Loreto Andrea Maldonado Mamut, Jens Dieter Kjær Modvig, Victoria Zillmer Payne, Anne Agathe Pedersen, Laura López Acedo, Agne Suminaite, Mahammed Zeeshan Siddique, Renjue Sun, Søren Winkel Holm, Peter Øllgaard Vilhelmsen, Benjamin Søndberg, Christian Noes-Rasmussen, Kabir Khanna, Samia Siddique Sama, Yiming Sun, ZhengZhong Sun, Miklós Kristóf Jásdi, Laia Poqui i Sallés, Nilas Tim Schüsler, Olivier Stephane Bonde Petersen, Jacob Kæstel-Hansen, Andreas Alsing Hauschild, Anders Thuelund Jensen, Simon Busch Iversen, Martin Nemes, Marius Spagl, Isaac Irani, Miguel Temboury Gutierrez, Kristian Sofus Knudsen, Carina Thusgaard Refsgaard, Martin Enggaard Kristensen, Anton Juhl, Zhaofeng Cai, François Gruwé, Michael Koch, Nicolaj Hostrup Langkjær, Lars Lohmann, Oldouz Majidi, Miriam Mazzeo

Reading material: Chapter 14, Chapter 15



Lecture Schedule

① Introduction

3 September: C1

Data: Feature extraction, and visualization

② Data, feature extraction and PCA

10 September: C2, C3

③ Measures of similarity, summary statistics and probabilities

17 September: C4, C5

④ Probability densities and data Visualization

24 September: C6, C7

Supervised learning: Classification and regression

⑤ Decision trees and linear regression

1 October: C8, C9 (Project 1 due before 13:00)

⑥ Overfitting, cross-validation and Nearest Neighbor

8 October: C10, C12

⑦ Performance evaluation, Bayes, and Naive Bayes

22 October: C11, C13

Piazza online help: <https://piazza.com/dtu.dk/fall2019/02450>

Videos/streaming of lectures: <https://video.dtu.dk>

⑧ Artificial Neural Networks and Bias/Variance

29 October: C14, C15

⑨ AUC and ensemble methods

5 November: C16, C17

Unsupervised learning: Clustering and density estimation

⑩ K-means and hierarchical clustering

12 November: C18 (Project 2 due before 13:00)

⑪ Mixture models and density estimation

19 November: C19, C20

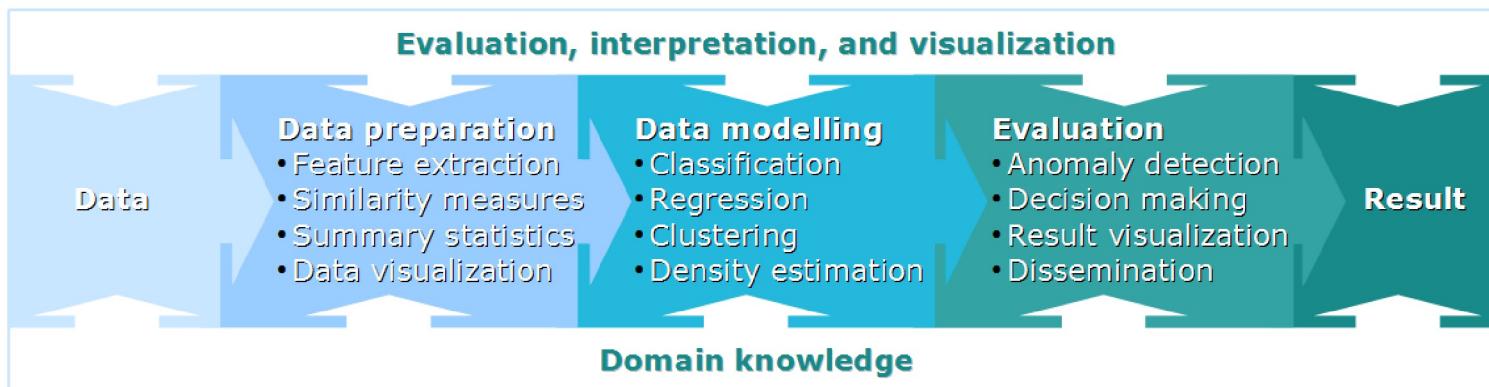
⑫ Association mining

26 November: C21

Recap

⑬ Recap and discussion of the exam

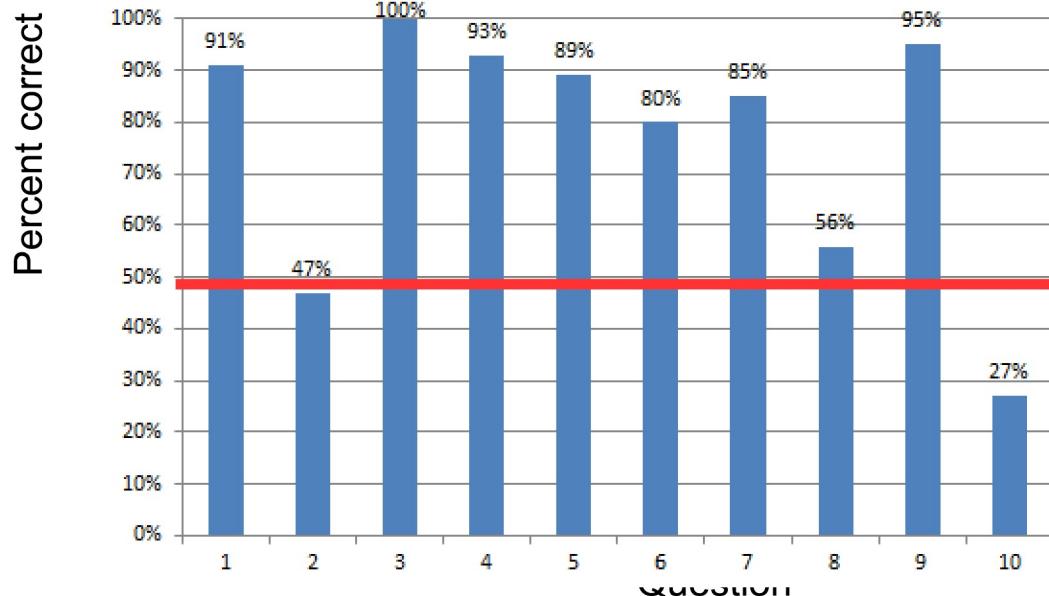
3 December: C1-C21 (Project 3 due before 13:00)



Learning Objectives

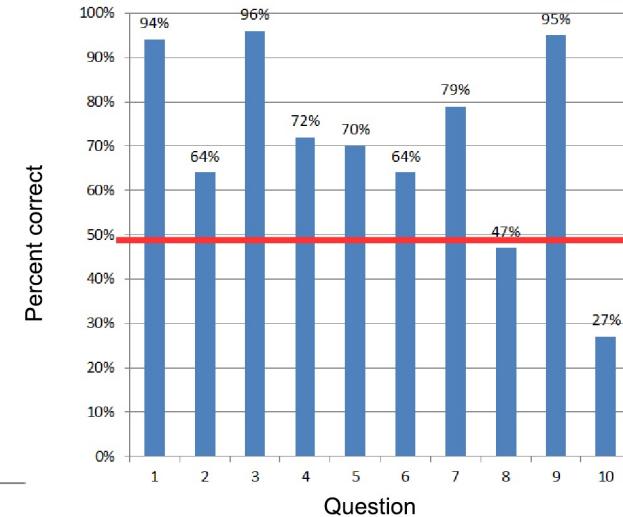
- Understand the Bias-Variance decomposition
- Understand and apply regularized least squares regression (i.e. ridge regression)
- Understand the principles behind artificial neural networks (ANNs) and how ANNs can be used for classification and regression
- Understand how logistic regression and ANNs can be extended to multi-class classification

Midterm practice test results

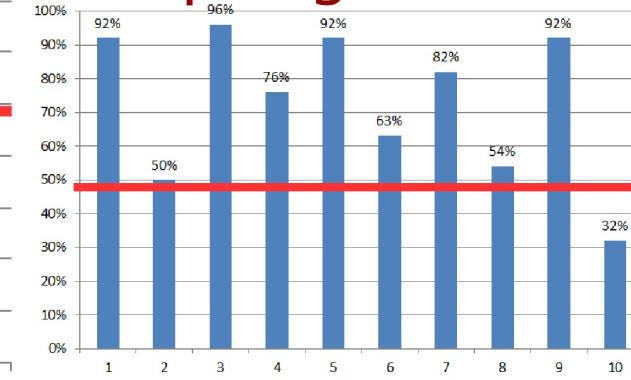


5 DTU Compute

DTU
Fall 2018



Spring 2019

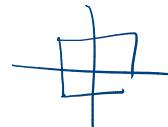


Lecture 8 29 October, 2019

Question 2:

Consider the classification problem given in figure 1 and the Decision Tree in figure 2 with two decisions denoted A and B. We will let \mathbf{x}_n define the x_1 and x_2 coordinates of a given observation whereas $\mathbf{x}_n - 0.5 \cdot \mathbf{1}$ denotes the subtraction of 0.5 from x_1 and x_2 .

Which one of the following classification rules would lead to a correct classification of the data?



- A: A: $\|\mathbf{x}_n - 0.5 \cdot \mathbf{1}\|_1 \leq 0.25$, B: $\|\mathbf{x}_n\|_\infty \leq 1$
- B: A: $\|\mathbf{x}_n\|_1 \leq 1$, B: $\|\mathbf{x}_n - 0.5 \cdot \mathbf{1}\|_2 \leq \infty$
- C: A: $\|\mathbf{x}_n - 0.5 \cdot \mathbf{1}\|_2 \leq 0.25$, B: $\|\mathbf{x}_n\|_\infty \leq 1$
- D: A: $\|\mathbf{x}_n - 0.5 \cdot \mathbf{1}\|_\infty \leq 0.25$, B: $\|\mathbf{x}_n\|_1 \leq 1$
- E: Don't know

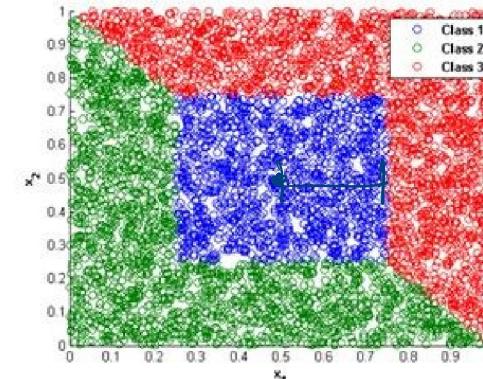


Figure 1

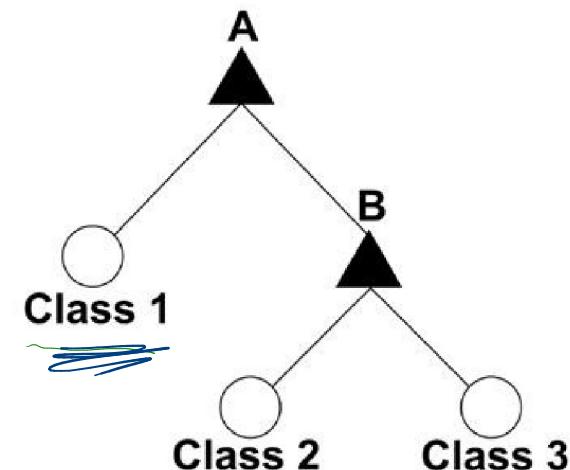


Figure 2
Lecture 8 29 October, 2019

$$\tilde{X} = U \Sigma V^T$$
$$\Sigma = \begin{bmatrix} s_1 & & & \\ & s_2 & & \\ & & s_3 & \\ & & & s_4 \end{bmatrix}$$

Question 8:

When carrying out a principal component analysis of a dataset with four attributes we obtain the following singular values $s_1=4$, $s_2=2$, $s_3=1$, and $s_4=0$. Which one of the following statements is wrong?

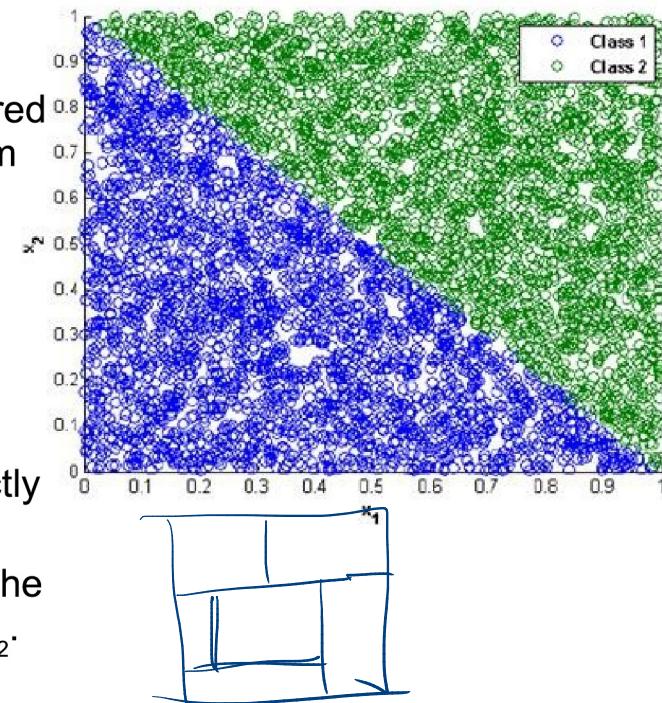
- A: The first principal component accounts for more than 60 % of the variation in the data.
- B: The third principal component accounts for less than 5 % of the variation in the data.
- C: The second principal component accounts for more than 20 % of the variation in the data.
- D: The data can be perfectly represented in a three dimensional sub-space.
- E: Don't know.

$$\frac{s_2^2}{s_1^2 + \dots + s_4^2}$$

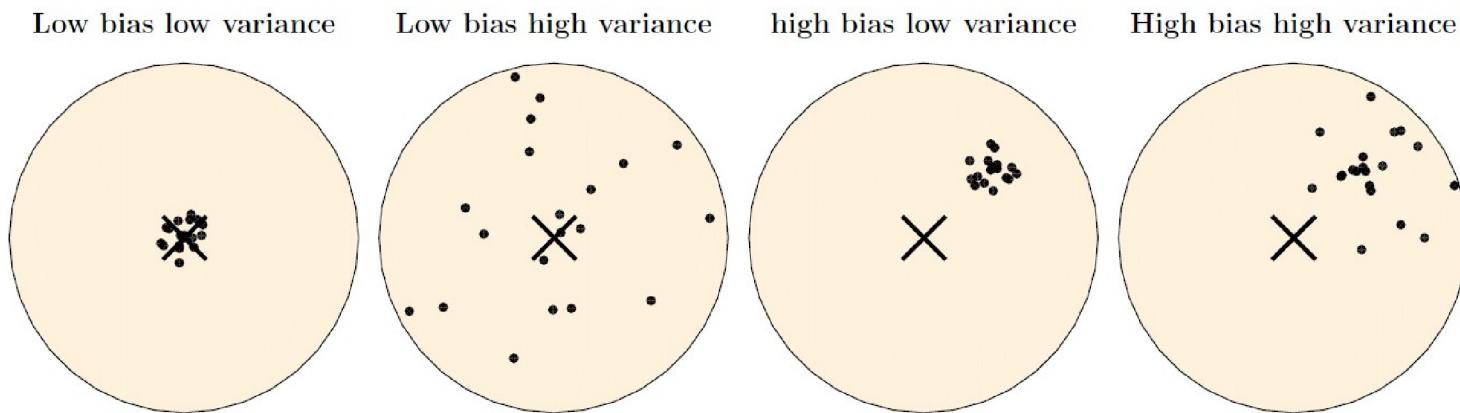
Question 10:

Consider the classification problem given in Figure 5 where x_1 and x_2 are used as features for a logistic regression classifier and a decision tree. The considered logistic regression models all include the constant term w_0 . Which one of the following statements is wrong?

- A: The two classes can be perfectly separated by a logistic regression model using x_1 and x_2 as features.
- B: A decision tree with less than five nodes can perfectly separate the classes using only x_1 and x_2 as features.
- C: A logistic regression model can perfectly separate the two classes using only the feature t given by $t = x_1 + x_2$.
- D: In logistic regression the probability that each observation belong to the two classes can be derived from the logistic function.
- E: Don't know.



What is bias and what is variance?



Regularized least squares

$$\mathbf{y} = \tilde{\mathbf{x}}^\top \mathbf{w} \quad \mathbf{w} = (\tilde{\mathbf{x}}^\top \tilde{\mathbf{x}})^{-1} \tilde{\mathbf{x}}^\top \mathbf{y}.$$

- Recall cost function from linear regression

$$E(\mathbf{w}) = \|\mathbf{y} - \tilde{\mathbf{X}}\mathbf{w}\|^2$$

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\mathbf{w} = \underbrace{\begin{bmatrix} \mathbf{I} \\ \vdots \\ \mathbf{I} \end{bmatrix}}_M \begin{bmatrix} \mathbf{w} \end{bmatrix}$$

- A parsimonious model can be obtained by **forcing** parameters towards zero.
- Problem: Columns of \mathbf{X} have very different scale (i.e. require large/small values of \mathbf{w})
- Therefore, standardize \mathbf{X} :

$$\hat{X}_{ij} = \frac{X_{ij} - \mu_j}{\hat{s}_j}, \quad \mu_j = \frac{1}{N} \sum_{i=1}^N X_{kj}, \quad \hat{s}_j = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (X_{ij} - \mu_j)^2}$$

- Note $\hat{\mathbf{X}}$ contains no constant term.

- Introduce regularization term $\lambda\|\mathbf{w}\|^2$ to penalize large weights:

$$E_\lambda(\mathbf{w}, w_0) = \sum_{i=1}^N (y_i - w_0 - \hat{\mathbf{x}}^\top \mathbf{w})^2 + \lambda\|\mathbf{w}\|^2 = \left\| \mathbf{y} - w_0 \mathbf{1} - \hat{\mathbf{X}} \mathbf{w} \right\|^2 + \lambda\|\mathbf{w}\|^2$$

- We can solve for w_0 and \mathbf{w} :

$$\frac{dE_\lambda}{dw_0} = \cancel{N} \sum_{i=1}^N -2(y_i - w_0 - \hat{\mathbf{x}}_i^\top \mathbf{w}) = -2N\mathbb{E}[y] - 2Nw_0 - N \left(\frac{1}{N} \sum_{i=1}^N \hat{\mathbf{x}}_i^\top \right) \mathbf{w} = 0$$

$$\Rightarrow w_0 = \mathbb{E}[y]$$

- With $\hat{y}_i = y_i - \mathbb{E}[y]$

$$E_\lambda = \left\| \hat{\mathbf{y}} - \hat{\mathbf{X}} \mathbf{w} \right\|^2 + \lambda\|\mathbf{w}\|^2$$

$$\begin{aligned} \frac{\partial E_\lambda}{\partial \mathbf{w}} &= \hat{\mathbf{X}}^\top (2(\hat{\mathbf{y}} - \hat{\mathbf{X}} \mathbf{w}) + 2\lambda \mathbf{w}) = 0 \\ &= 2(\hat{\mathbf{X}}^\top \hat{\mathbf{y}} + 2\lambda \mathbf{I}) \mathbf{w} = 2\hat{\mathbf{X}}^\top \hat{\mathbf{y}} \\ \mathbf{w} &= \underbrace{(\hat{\mathbf{X}}^\top \hat{\mathbf{X}} + \lambda \mathbf{I})^{-1}}_{\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}} \hat{\mathbf{X}}^\top \hat{\mathbf{y}}. \end{aligned}$$

- Setting the derivative wrt. \mathbf{w} equal to zero and solving for \mathbf{w} yields

$$\mathbf{w}^* = (\hat{\mathbf{X}}^\top \hat{\mathbf{X}} + \lambda \mathbf{I}) \setminus (\hat{\mathbf{X}}^\top \hat{\mathbf{y}})$$

Selecting λ

$$\begin{aligned} VV^T &= I \\ UU^T &= U^TU = I \end{aligned}$$



$$(AB)^{-1} = B^{-1}A^{-1}$$

$$\hat{X} = U\Sigma V^T$$

$$\begin{aligned} w &= ((U\Sigma V^T)^T(U\Sigma V^T) + \lambda I)^{-1}(U\Sigma V^T)^T y \\ &= (V\Sigma^T U^T U\Sigma V^T + VV^T)^{-1}(V\Sigma V^T)^T y \\ &= (V(\begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ 0 & & \sigma_M^2 \end{bmatrix} + \lambda I)V^T)^{-1}(V\Sigma V^T)^T y \end{aligned}$$

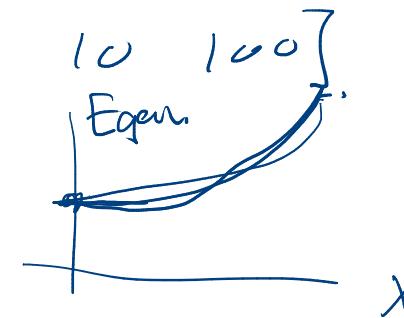
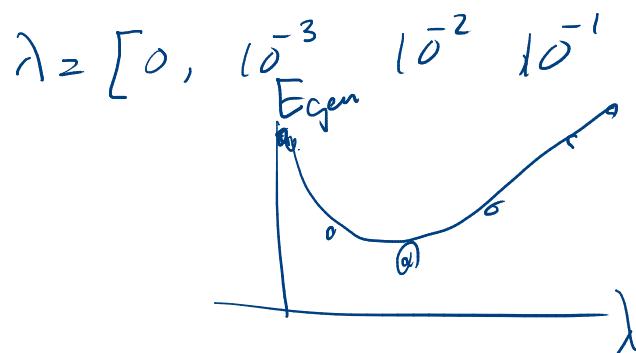
$$w^* = (\hat{X}^T \hat{X} + \lambda I) \backslash (\hat{X}^T \hat{y}) \propto \frac{Xy}{X^2 + \lambda} = \left(V \begin{bmatrix} \frac{1}{\sigma_1^2 + \lambda} & & 0 \\ 0 & \ddots & \\ & & \frac{1}{\sigma_M^2 + \lambda} \end{bmatrix} V^T \right) (V\Sigma^T Vy)$$

- Suppose

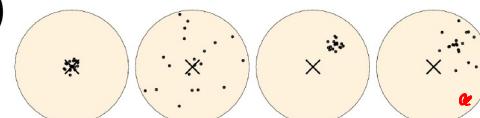
So if $\lambda = 0$ then no effect, else if $\lambda \rightarrow \infty$ then $w^* \rightarrow 0$

λ controls complexity of model. Select λ using cross-validation

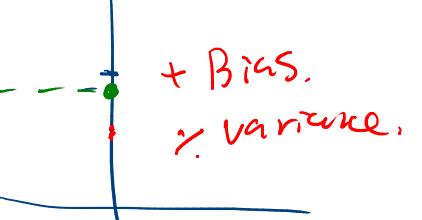
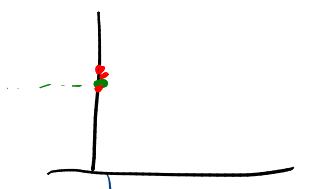
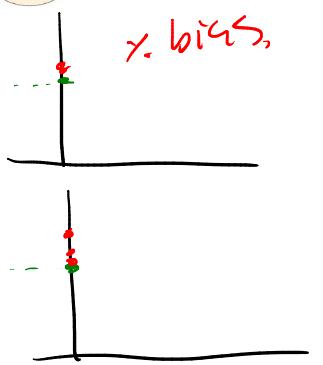
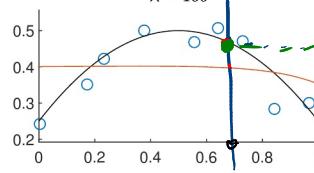
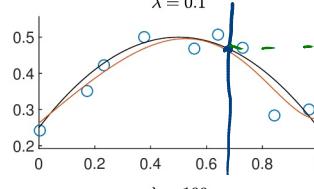
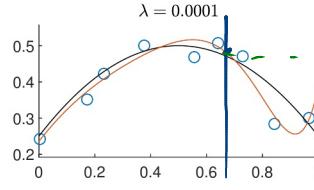
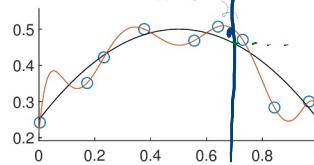
$$= V \left(\begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & 0 \\ 0 & \ddots & \\ & & \frac{\sigma_M^2}{\sigma_M^2 + \lambda} \end{bmatrix} \right) \circ (Vy)$$



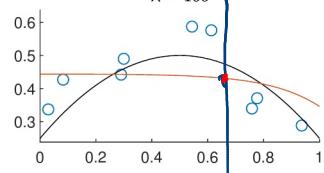
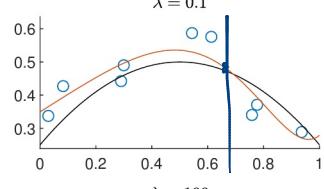
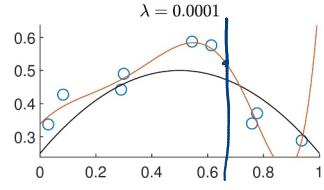
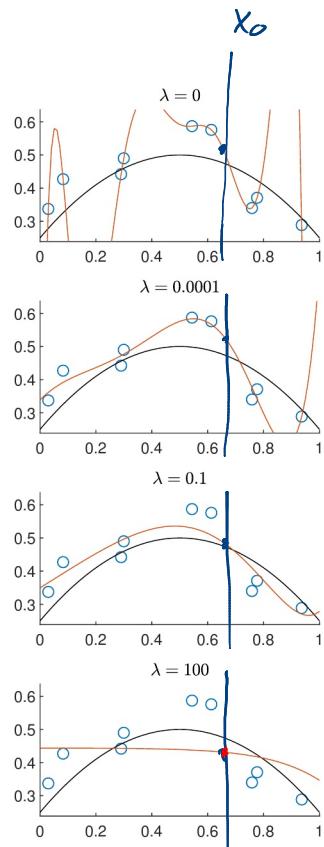
How does different values of λ (vertical) affect the bias/variance of learned function (red lines)



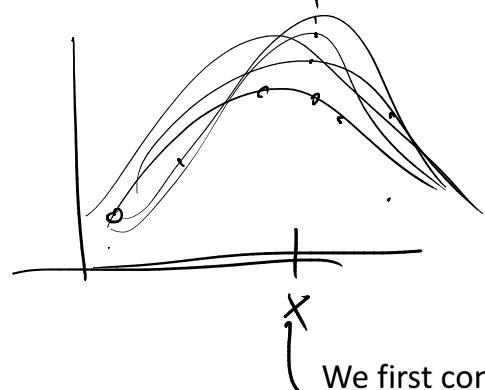
\downarrow variance +
 \uparrow bias,



+ Bias.
 γ variance.



The Bias-Variance decomposition



We first consider x fixed

$$\mathbb{E}_{\mathcal{D},y|x} [(y - f_{\mathcal{D}}(x))^2]$$

$$= \mathbb{E}_{\mathcal{D},y|x} [(y - \bar{y}(x) + \bar{y}(x) - f_{\mathcal{D}}(x))^2] = \mathbb{E} \left[(y - \bar{g}(x))^2 + (\bar{g}(x) - f_{\mathcal{D}}(x))^2 + (y - \bar{g}(x))(\bar{g}(x) - f_{\mathcal{D}}(x)) \right]$$

$$= \mathbb{E}_{y|x} [(y - \bar{y}(x))^2] + \mathbb{E}_{\mathcal{D}} [(\bar{y}(x) - f_{\mathcal{D}}(x))^2] + 2 \cancel{\mathbb{E}_{\mathcal{D},y|x} [(\bar{y}(x) - f_{\mathcal{D}}(x))(\bar{y}(x) - f_{\mathcal{D}}(x))]}$$

$$\mathbb{E}_x [f(x) + g(x)] = \mathbb{E}[f] + \mathbb{E}[g]$$

$$\mathbb{E}_{z|x} [g(x)f(z)] = \mathbb{E}_z [\mathbb{E}_x [g(x)f(z)]] = \mathbb{E}_z [f(z)\mathbb{E}_x [g(x)]] \\ \geq \mathbb{E}_z [f_z]\mathbb{E}_x [g(x)]$$

$$= \int_{\mathcal{D}} p(y|x) \int_{\mathcal{D}} p(x,y) (y - f_{\mathcal{D}}(x))^2 dx dy$$

T training sets, \mathcal{D}_t .
N test observations $\mathcal{D}_t^{\text{test}}$.
train model $f_{\mathcal{D}_t}$ on \mathcal{D}_t

$$\mathbb{E}_{\mathcal{D}} [E^{\text{gen}}] = \mathbb{E}_{\mathcal{D},(x,y)} [(y - f_{\mathcal{D}}(x))^2]$$

$$= \frac{1}{T} \sum_{t=1}^T (E_{\mathcal{D}_t^{\text{test}}, f_{\mathcal{D}_t}})$$

$$= \frac{1}{T} \sum_{t=1}^T \left(\frac{1}{N} \sum_{i=1}^N (y_i^{\text{t,test}} - \bar{f}_{\mathcal{D}_t}(x_i^{\text{t,test}}))^2 \right)$$

$$\bar{y}(x) = \underline{\mathbb{E}_{y|x} [y]}$$

$$\mathbb{E}_{\mathcal{D},y|x} [(y - \bar{y}(x))^2] = \int p(g|x) y dy$$

$$\mathbb{E}[(y - \bar{g}(x))^2] = \mathbb{E}[y^2] - \bar{g}(x)^2 = \bar{g}(x) - \bar{g}(x) = 0$$



The Bias-Variance decomposition



$$\mathbb{E}_{\mathcal{D}} [E^{\text{gen}}] = \mathbb{E}_{\mathcal{D},(\mathbf{x},y)} \left[(y - f_{\mathcal{D}}(\mathbf{x}))^2 \right]$$

We first consider \mathbf{x} fixed

$$\begin{aligned} & \mathbb{E}_{\mathcal{D},y|\mathbf{x}} \left[(y - f_{\mathcal{D}}(\mathbf{x}))^2 \right] & \bar{y}(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}} [y] \\ &= \mathbb{E}_{\mathcal{D},y|\mathbf{x}} \left[(y - \bar{y}(\mathbf{x}) + \bar{y}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E}_{y|\mathbf{x}} \left[(y - \bar{y}(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right] + \cancel{2\mathbb{E}_{\mathcal{D},y|\mathbf{x}} [(y - \bar{y}(\mathbf{x})) (\bar{y}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))]} \end{aligned}$$



The Bias-Variance decomposition



$$\mathbb{E}_{\mathcal{D}, y|\mathbf{x}} \left[(y - f_{\mathcal{D}}(\mathbf{x}))^2 \right] = \mathbb{E}_{y|\mathbf{x}} \left[(y - \bar{y}(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right]$$

$$\begin{aligned} &\rightarrow \mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}) + \bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathcal{D}} \left[(\bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right] + 2\mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x})) (\bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x})) \right] \end{aligned}$$

$$\underbrace{\bar{f}(\mathbf{x})}_{\hat{f}(\mathbf{x})} = \mathbb{E}_{\mathcal{D}} [f_{\mathcal{D}}(\mathbf{x})]$$
$$\underbrace{\bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x})}_{\hat{f}(\mathbf{x}) - \bar{f}(\mathbf{x})} = \underbrace{\bar{f}(\mathbf{x}) - \bar{f}(\mathbf{x})}_{0} = 0$$



The Bias-Variance decomposition



$$\mathbb{E}_{\mathcal{D}, y|\mathbf{x}} \left[(y - f_{\mathcal{D}}(\mathbf{x}))^2 \right] = \mathbb{E}_{y|\mathbf{x}} \left[(y - \bar{y}(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right]$$

$$\text{Var}_{\mathbf{x}} [f] = \mathbb{E}_{\mathbf{x}} \left[(f(\mathbf{x}) - \mathbb{E}[f])^2 \right]$$

$$\bar{f}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} [f_{\mathcal{D}}(\mathbf{x})]$$

$$\mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right]$$

$$= \mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}) + \bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right]$$
~~$$= \cancel{\mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 \right]} + \mathbb{E}_{\mathcal{D}} \left[(\bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right] + 2\mathbb{E}_{\mathcal{D}} \left[(\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x})) (\bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x})) \right]$$~~

$$\mathbb{E}_{\mathcal{D}, y|\mathbf{x}} \left[(y - f_{\mathcal{D}}(\mathbf{x}))^2 \right]$$

$$= \mathbb{E}_{y|\mathbf{x}} \left[(y - \bar{y}(\mathbf{x}))^2 \right] + (\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 + \mathbb{E}_{\mathcal{D}} \left[(\bar{f}(\mathbf{x}) - f_{\mathcal{D}}(\mathbf{x}))^2 \right]$$

$$= \text{Var}_{y|\mathbf{x}} [y] + (\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 + \text{Var}_{\mathcal{D}} [f_{\mathcal{D}}(\mathbf{x})]$$

The Bias-Variance decomposition



$$p(x,y) \geq p(x)p(y|x)$$

$$\begin{aligned} \mathbb{E}_{x,y}[f(x,y)] &= \iint p(x,y) f(x,y) dx dy \\ &= \iint p(x) f(x,y) p(y|x) dy dx \\ &\geq \int p(x) \left(\int p(y|x) f(x,y) dy \right) dx = \mathbb{E}_x \left[\mathbb{E}_{y|x}[f(x,y)] \right] \end{aligned}$$

$$\mathbb{E}_{\mathcal{D}} [E^{\text{gen}}] = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{D}, y|\mathbf{x}} \left[(y - f_{\mathcal{D}}(\mathbf{x}))^2 \right] \right]$$

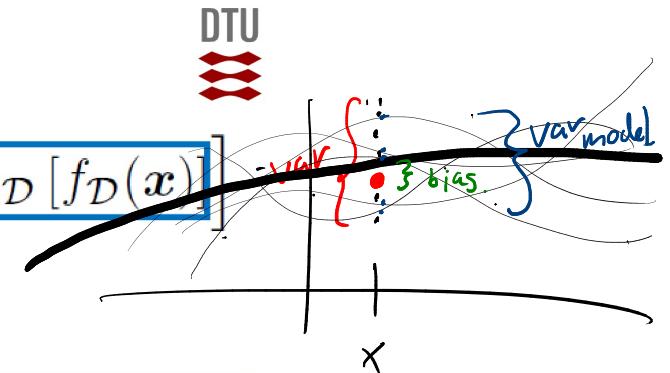
$$\mathbb{E}_{\mathcal{D}} [E^{\text{gen}}] = \mathbb{E}_{\mathbf{x}} \left[\text{Var}_{y|\mathbf{x}} [y] + (\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 + \text{Var}_{\mathcal{D}} [f_{\mathcal{D}}(\mathbf{x})] \right]$$



[https://commons.wikimedia.org/wiki/File:Frustrated_man_at_a_desk_\(cropped\).jpg](https://commons.wikimedia.org/wiki/File:Frustrated_man_at_a_desk_(cropped).jpg)

The Bias-Variance decomposition

$$\mathbb{E}_{\mathcal{D}} [E^{\text{gen}}] = \mathbb{E}_{\mathbf{x}} \left[\text{Var}_{y|\mathbf{x}} [y] + (\bar{y}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 + \text{Var}_{\mathcal{D}} [f_{\mathcal{D}}(\mathbf{x})] \right]$$

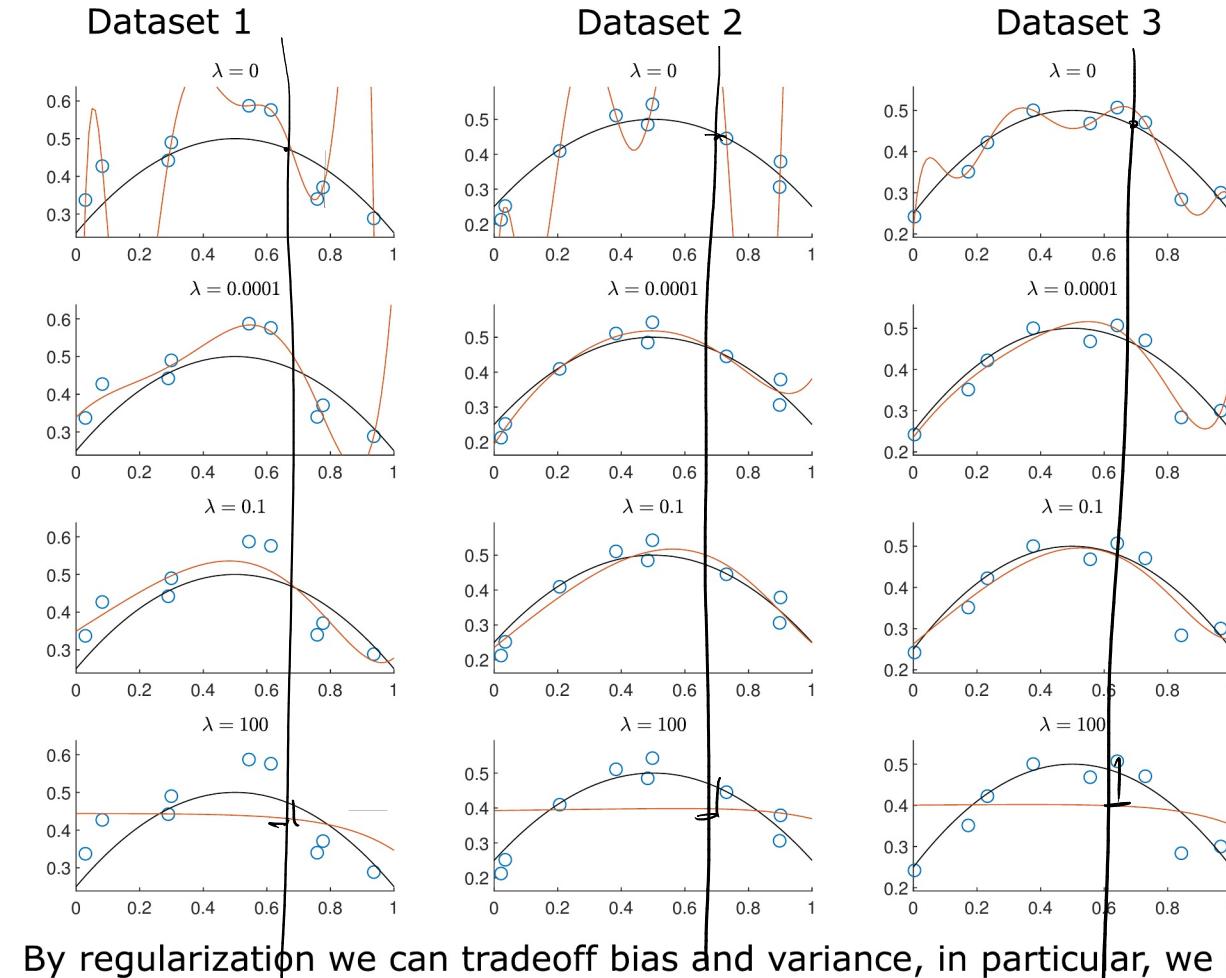


The first term does not depend at all upon our choice of model but simply represents the intrinsic difficulty of the problem. We cannot make this term any larger or smaller by selecting one model over another.

The second term is the **bias** term. It tells us how much the average values of models trained on different training datasets differ compared to the true mean of the data.

The third term is the **variance** term. It tells us how much the model wiggles when trained on different sets of training data. That is, when you train the models on N different (random) sets of training data and the models (the prediction curves) are nearly the same this term is small.

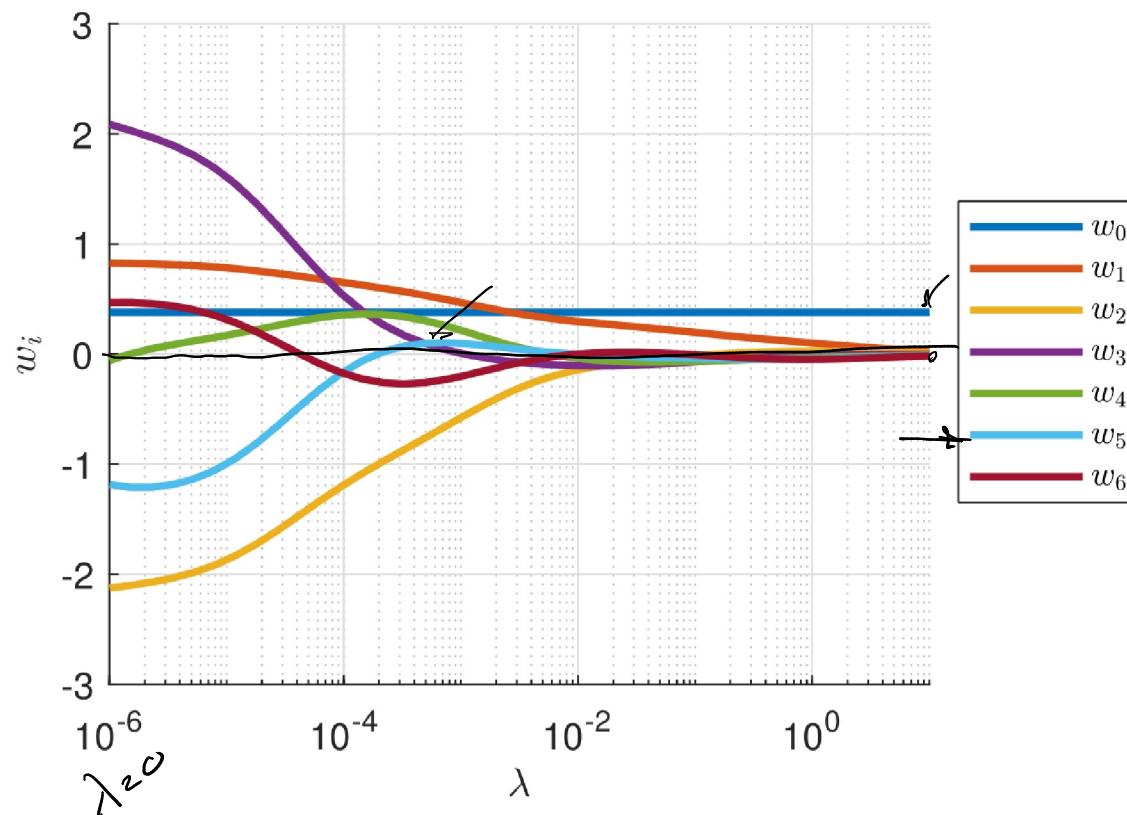
The bias variance decomposition



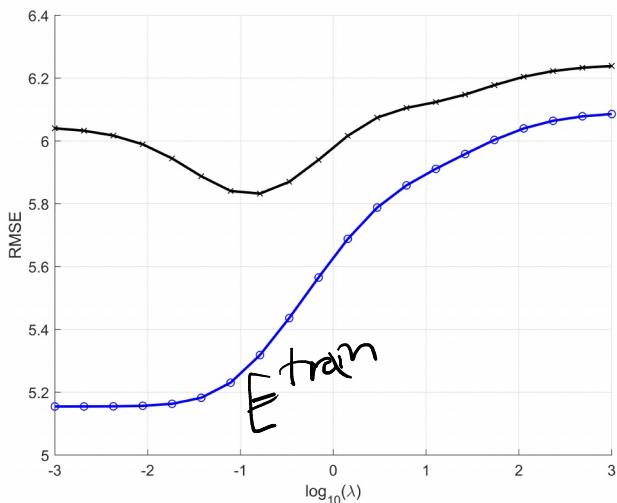
By regularization we can tradeoff bias and variance, in particular, we can hope to substantially reduce variance without introducing too much bias!

Parameters w^* as function of λ

$$E_\lambda(\mathbf{w}) = \sum_{i=1}^N (\hat{y}_i - w_0 - \hat{\mathbf{x}}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|^2$$



Quiz 1, Bias-variance (Fall 2017)



Using 54 observations of a dataset about Basketball, we would like to predict the average points scored per game (y) based on the four features. For this purpose we consider regularized least squares regression which minimizes with respect to \mathbf{w} the following cost function:

$$E(\mathbf{w}) = \sum_n (y_n - [1 \ x_{n1} \ x_{n2} \ x_{n3} \ x_{n4}] \mathbf{w})^2 + \lambda \mathbf{w}^\top \mathbf{w},$$

We consider 20 different values of λ and use leave-

one-out cross-validation to estimate the performance (measured by mean-squared error) of each of these different values of λ and plot the result in the figure. For the value of $\lambda = 0.6952$ the following model is identified:

$$f(\mathbf{x}) = 2.76 - 0.37x_1 + 0.01x_2 + 7.67x_3 + 7.67x_4.$$

Which one of the following statements is correct?

- A. In the figure the blue curve with circles corresponds to the training error whereas the black curve with crosses corresponds to the test error.
- B. According to the model defined for $\lambda = 0.6952$ increasing a players height x_1 will increase his average points scored per game.
- C. There is no optimal way of choosing λ since increasing λ reduces the variance but increases the bias.
- D. As we increase λ the 2-norm of the weight vector \mathbf{w} will also increase.
- E. Don't know.

The correct answer is *A*: The blue curve monotonically increases with λ reflecting a worse fit to the training set as we increase λ using regularization we can reduce the variance by introducing bias and the black curve indicates that an optimal tradeoff at around $10^{-0.8}$ as reflected by the test error indicated in the black curve being minimal. As we increase λ we will

penalize the weights according to the squared 2-norm more and more and thus the 2-norm will be reduced. Finally, according to the fitted model we observe that the coefficient in front of x_1 (Height) is negative thus indicating that an increase in height will reduce the models prediction of average points scored per game.

General linear model

DTU

$$g_{\text{link}}(z) = z, \quad d(y, \hat{y}) = \|y - \hat{y}\|^2$$

$$g_{\text{link}}(z) = \sigma(z), \quad d(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

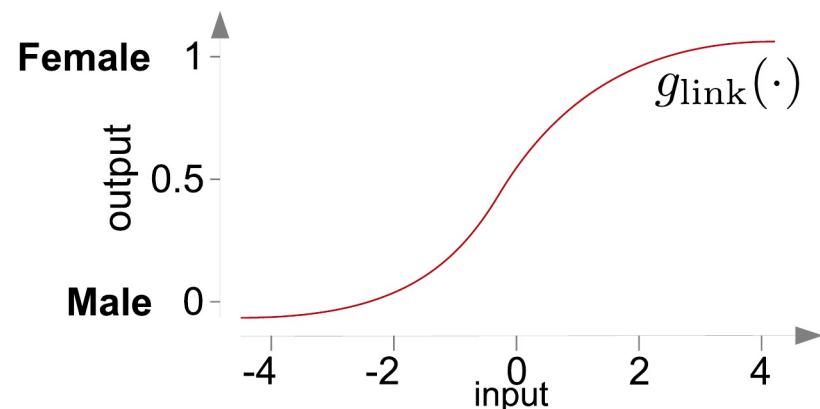
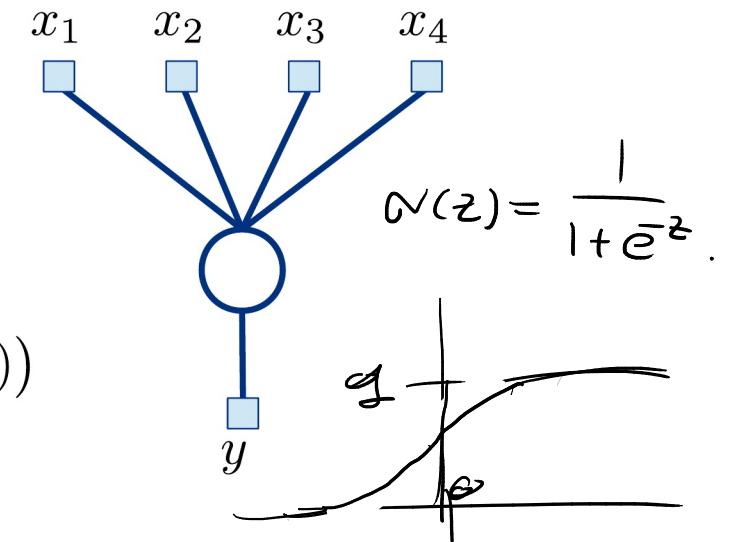
- Remember the generalized linear model?

- Data $\{\mathbf{x}_n, y_n\}_{n=1}^N$

- Model $f(\mathbf{x}) = g_{\text{link}}(\mathbf{x}^\top \mathbf{w})$

- Cost function $d(y, f(\mathbf{x}))$

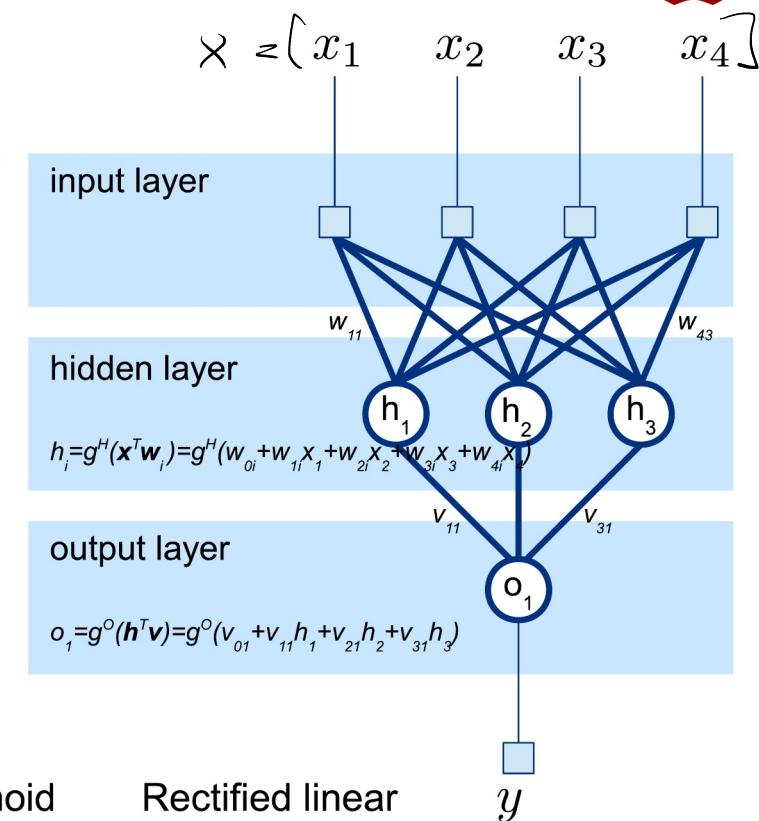
- Parameters $\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{n=1}^N d(y_n, f(\mathbf{x}_n))$



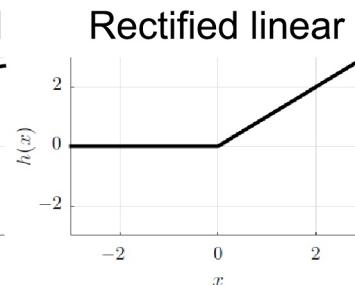
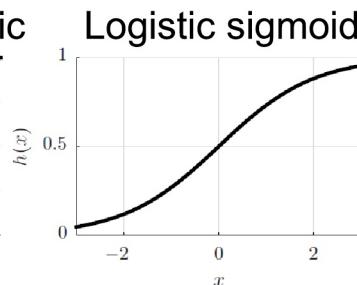
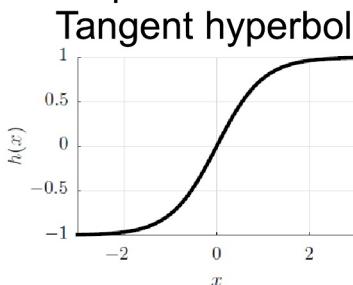
Artificial neural networks

Feed forward network

- Each “neuron”
 - Computes a non-linear function of the sum of its inputs
 - Is just like a generalized linear model
 - Has its own set of parameters
- Modeling choices
 - Cost function
 - Non-linearities
 - Number of neurons and hidden layers
 - Selection of inputs
- Parameter estimation using numerical optimization methods
- Very flexible model: Can easily overfit



Example of non-linearities:



Data: $\{\mathbf{x}_i, y_i\}$

Model: $f(\mathbf{x}) = h^{(2)} \left(v_{10} + \sum_{j=1}^H v_{1j} h^{(1)} (\tilde{\mathbf{x}}^\top \mathbf{w}_j) \right)$

Distance: $d(y, f(\mathbf{x}))$

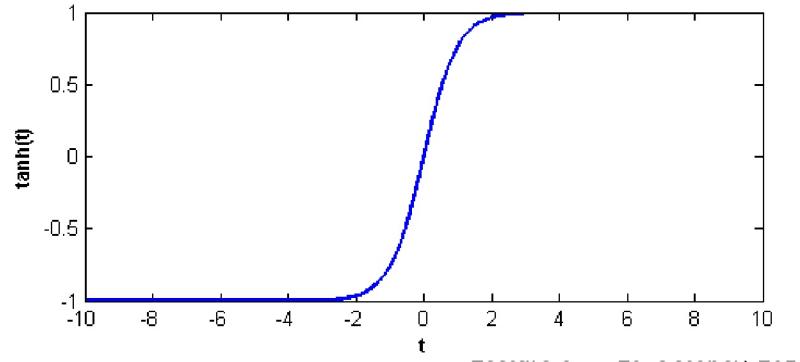
Cost: $E = \underbrace{\sum_{i=1}^N d(y_i, f(\mathbf{x}_i))}_{\text{Cost}}$

Common choices

$$h^{(1)}(x) = \tanh(x)$$

$$h^{(2)}(x) = x$$

$$d(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$$



Neurons and layers

Recall:

$$f(\mathbf{x}) = h^{(2)} \left(v_{10} + \sum_{j=1}^H v_{1j} h^{(1)} (\tilde{\mathbf{x}}^\top \mathbf{w}_j) \right)$$

- Let $z_j^{(1)}$ be output of j 'th hidden unit

$$z_j^{(1)} = h^{(1)} \left(\mathbf{w}_j^{(1)\top} \tilde{\mathbf{x}} \right)$$

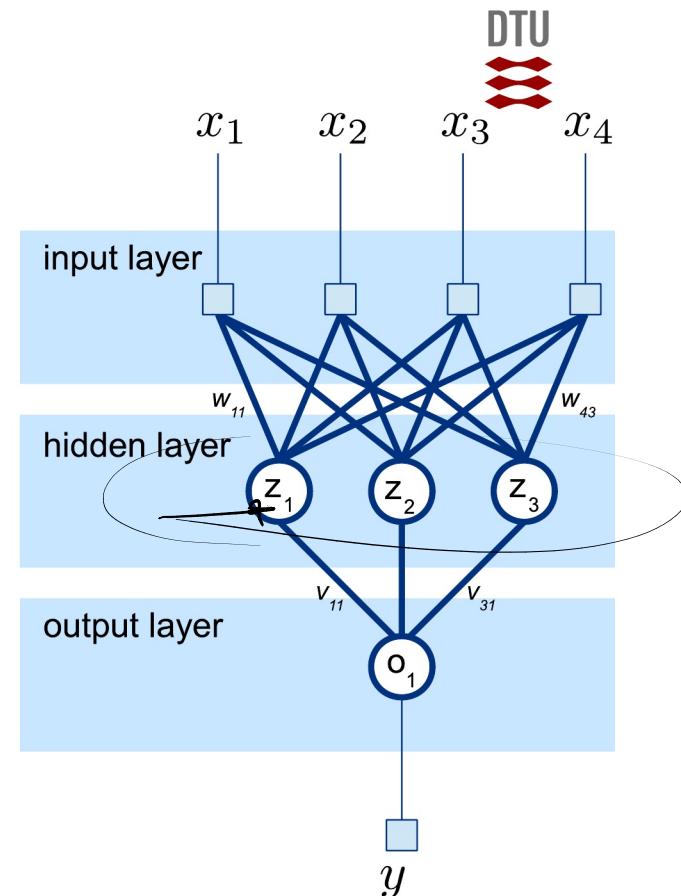
Abbreviated $z^{(1)} = h^{(1)} \left(\mathbf{W}^{(1)} \tilde{\mathbf{x}} \right)$

- Output $\begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ \vdots \\ z_H^{(1)} \end{bmatrix} = h^{(1)} \left(\begin{bmatrix} (\mathbf{w}_1^{(1)})^\top \\ \vdots \\ (\mathbf{w}_H^{(1)})^\top \end{bmatrix} \tilde{\mathbf{x}} \right)$

$$f(\mathbf{x}) = h^{(2)} \left(v_{10} + \sum_{j=1}^H v_{1j} z_j^{(1)} \right) = h^{(2)} \left(\mathbf{W}^{(2)} \tilde{\mathbf{z}}^{(1)} \right)$$

$$\mathbf{W}^{(2)} = [v_{10} \quad v_{11} \quad v_{12} \dots \quad v_{1H}]$$

We consider each $z_j^{(1)}$ a neuron and $z^{(1)}$ a (hidden) layer



Quiz 2, Artificial Neural Network (Fall 2017)

We will consider an artificial neural network (ANN) trained to predict the average score of a player (i.e., y). The ANN is based on the model:

$$f(\mathbf{x}, \mathbf{w}) = w_0^{(2)} + \sum_{j=1}^2 w_j^{(2)} h^{(1)}([1 \ \mathbf{x}] \mathbf{w}_j^{(1)}).$$

where $h^{(1)}(x) = \max(x, 0)$ is the rectified linear function used as activation function in the hidden layer (i.e., positive values are returned and negative values are set to zero). We will consider an ANN with two hidden units in the hidden layer defined by:

$$\mathbf{w}_1^{(1)} = \begin{bmatrix} 21.78 \\ -1.65 \\ 0 \\ -13.26 \\ -8.46 \end{bmatrix}, \quad \mathbf{w}_2^{(1)} = \begin{bmatrix} -9.60 \\ -0.44 \\ 0.01 \\ 14.54 \\ 9.50 \end{bmatrix},$$

and $w_0^{(2)} = 2.84$, $w_1^{(2)} = 3.25$, and $w_2^{(2)} = 3.46$.

What is the predicted average score of a basketball player with observation vector $\mathbf{x}^* = [6.8 \ 225 \ 0.44 \ 0.68]^\top$?

- A. 1.00
- B. 3.74
- C. 8.21
- D. 11.54
- E. Don't know.

The output is given by:

$$f(\mathbf{x}, \mathbf{w}) = 2.84$$

$$\begin{aligned}
 & + 3.25 \cdot \max([1 \ 6.8 \ 225 \ 0.44 \ 0.68]) \cdot \begin{bmatrix} 21.78 \\ -1.65 \\ 0 \\ -13.26 \\ -8.46 \end{bmatrix}, 0) \\
 & + 3.46 \max([1 \ 6.8 \ 225 \ 0.44 \ 0.68]) \cdot \begin{bmatrix} -9.60 \\ -0.44 \\ 0.01 \\ 14.54 \\ 9.50 \end{bmatrix}, 0) \\
 & = 2.84 + 3.25 \cdot \max(-1.027, 0) + 3.46 \max(2.516, 0) \\
 & = 11.54
 \end{aligned}$$

Generalization 1: Multiple outputs

- As before define: $z^{(1)} = h^{(1)}(\mathbf{W}^{(1)}\tilde{x})$

- Now let $\mathbf{W}^{(2)}$ be a $C \times H$ matrix then:

$$\mathbf{W}^{(2)} : 1 \times H \quad \mathbf{y} = \mathbf{f}(\mathbf{x}) = h^{(2)}(\mathbf{W}^{(2)}\tilde{z}^{(1)})$$

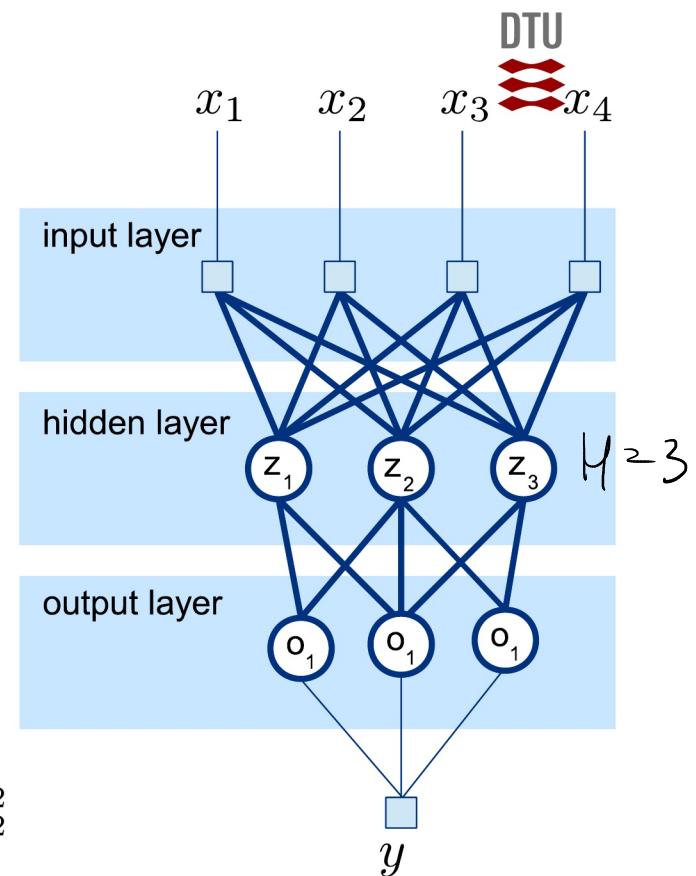
$\mathbf{y} : \left[\begin{array}{c} \vdots \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_C \end{array} \right]$

will be C -dimensional

- Re-define error function

$$E = \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i)\|_2^2$$

$$(\mathbf{x}, \mathbf{y}) \sim \left(\left[\begin{array}{c} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{array} \right], \left[\begin{array}{c} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_C \end{array} \right] \right)$$



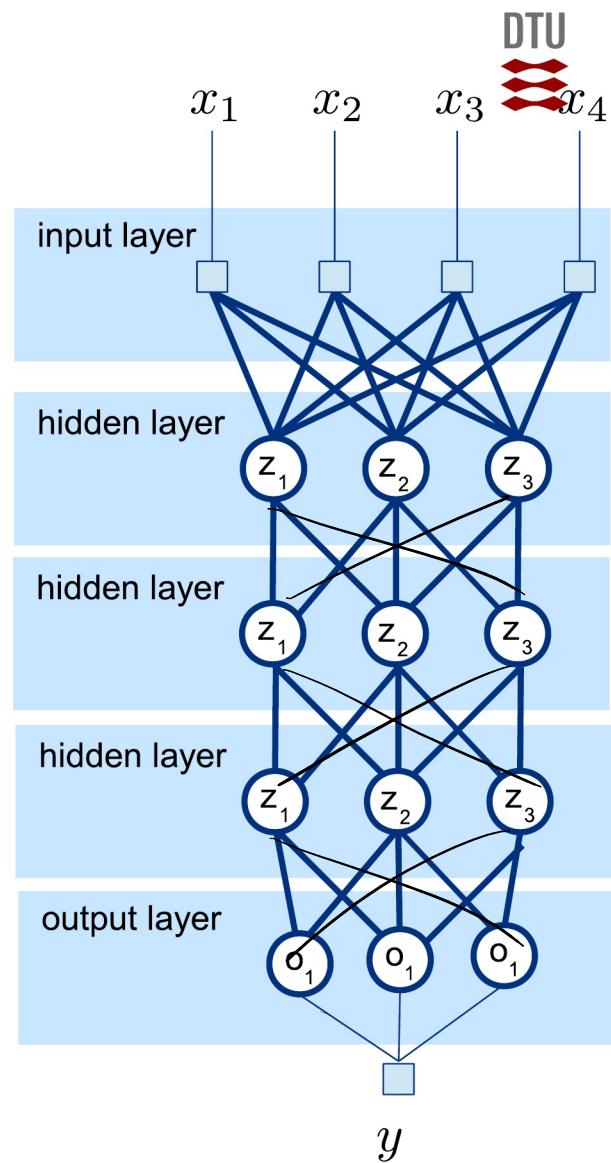
Generalization 2: Multiple layers

- Define $z^{(0)} = x$
- For each layer $l = 1, \dots, L$ compute

$$z_j^{(l)} = h^{(l)} \left(\mathbf{W}^{(l)} z^{(l-1)} \right)$$

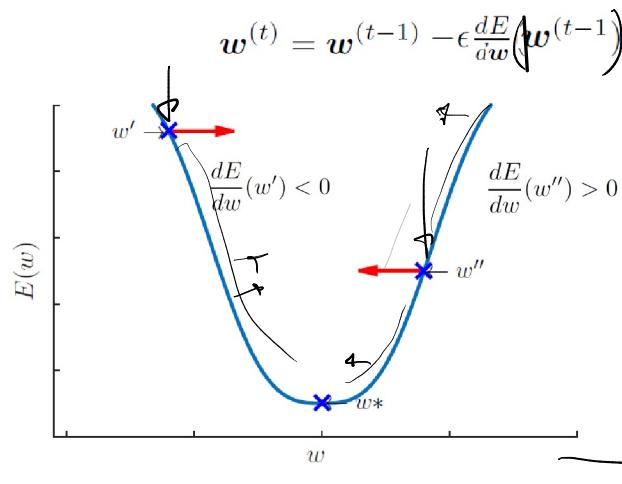
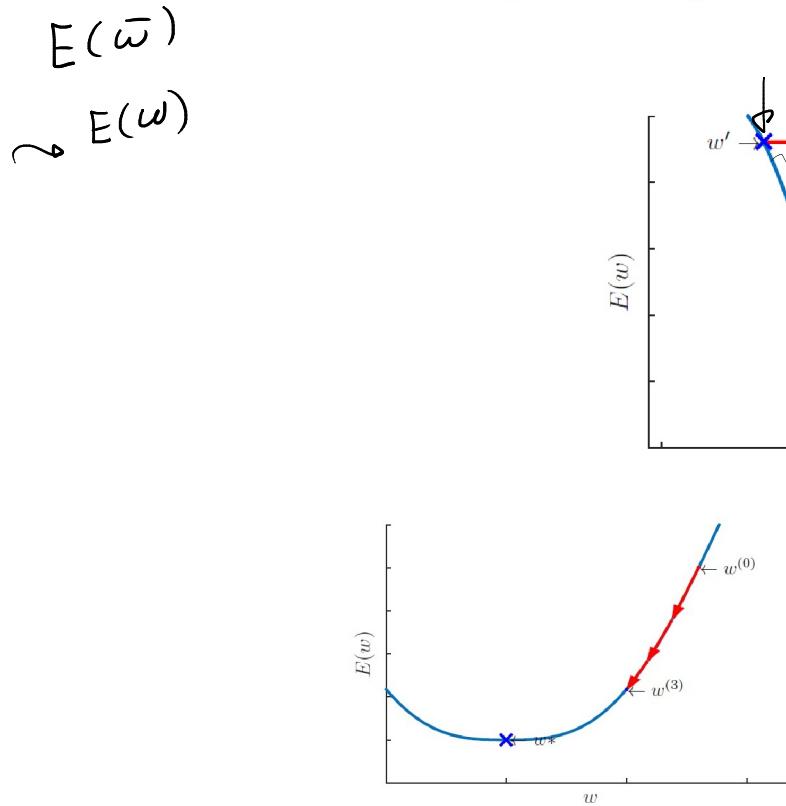
- Output is simply

$$f(x) = z^{(L)}$$



Gradient descent

- Start from an initial guess at $w^*, w^{(0)}$
- At step t of the algorithm, modify $w^{(t-1)}$ to produce a better guess $w^{(t)}$.



$$E(\bar{w}) = E(w^*) + [\nabla E(w^*)]^\top \times (\bar{w} - w^*)$$

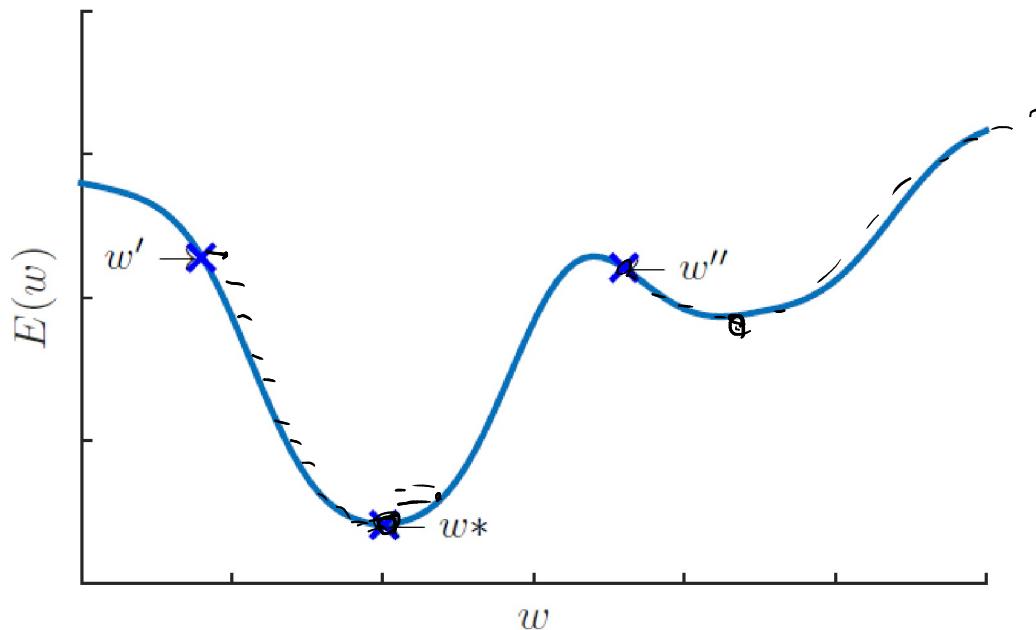
~~+ $\frac{1}{2} \nabla^2 E(x)$~~

$$E(\delta) = E(w^*) + \nabla E(w^*)^\top \delta.$$

guess. $\delta = -\epsilon \nabla E(w^*)$

$$\begin{aligned} & \xrightarrow{\sim} E(\delta = -\epsilon \nabla E(w^*)) \\ &= E(w^*) - \epsilon \|\nabla E(w^*)\| \end{aligned}$$

Contrary to least-squares linear regression and logistic regression ANNs have issues of local minima



Single and multi-class: One out of K coding

$$y = 1, \dots, k$$

Nationality

Denmark Norway Sweden

One-out-of-K coding

→

↗

TXT=		X_tmp=	Denmark	Norway	Sweden
			0	0	1
	'Sweden'	0	0	1	
	'Sweden'	0	0	1	
	'Sweden'	0	0	1	
	'Sweden'	0	0	1	
	'Norway'	0	1	0	
	'Norway'	0	1	0	
	'Norway'	0	1	0	
	'Norway'	0	1	0	
	'Sweden'	0	0	1	
	'Norway'	0	1	0	
	'Denmark'	1	0	0	
	'Denmark'	1	0	0	
	'Sweden'	0	0	1	
	'Sweden'	0	0	1	
	'Sweden'	0	0	1	
	'Denmark'	1	0	0	
	'Sweden'	0	0	1	
	'Norway'	0	1	0	
	'Denmark'	1	0	0	

Multi-class classification

- Logistic regression, $y = 0, 1$:

$$p(y|\theta) = \theta^y (1 - \theta)^{1-y} \quad \text{Ber}^n,$$

$$\theta = \sigma(\mathbf{x}^\top \mathbf{w}) \quad \theta \in [0, 1]$$

- Multinomial regression, $y = 1, 2, \dots, K$

z_k : one-of- K encoding of y ,

$$z_{k=1} = 1$$

$$p(y|\theta) = \prod_{i=1}^K \theta_k^{z_k}$$

categorical.

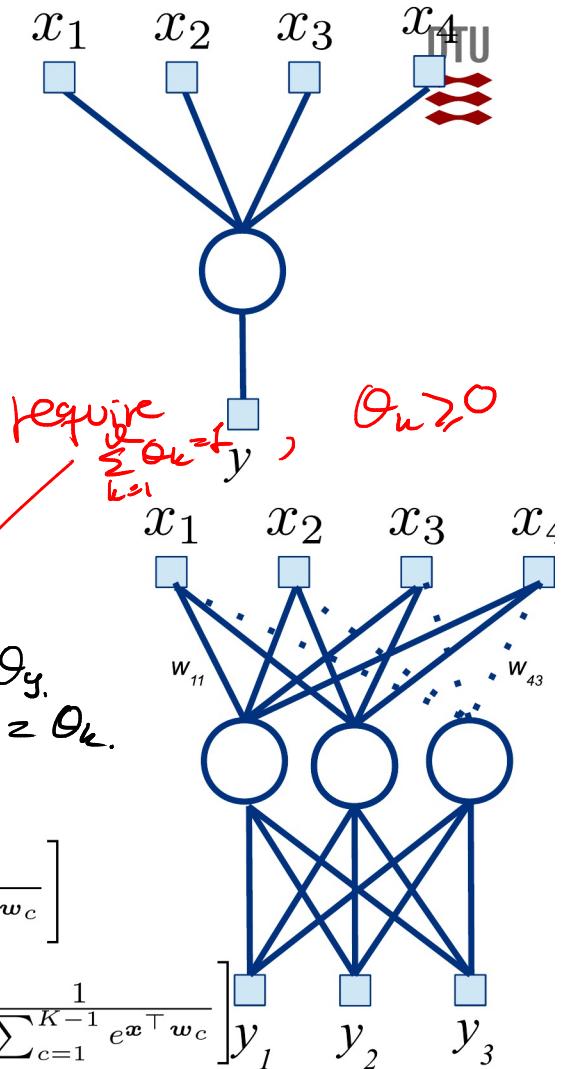
$$p(y|\theta) = \theta_y$$

$$p(y=k|\theta) = \theta_k$$

$$\theta = \text{softmax}([\mathbf{x}^\top \mathbf{w}_1 \quad \dots \quad \mathbf{x}^\top \mathbf{w}_K])$$

$$= \begin{bmatrix} \frac{e^{\mathbf{x}^\top \mathbf{w}_1}}{\sum_{c=1}^K e^{\mathbf{x}^\top \mathbf{w}_c}} & \dots & \frac{e^{\mathbf{x}^\top \mathbf{w}_{K-1}}}{\sum_{c=1}^K e^{\mathbf{x}^\top \mathbf{w}_c}} & \frac{e^{\mathbf{x}^\top \mathbf{w}_K}}{\sum_{c=1}^K e^{\mathbf{x}^\top \mathbf{w}_c}} \end{bmatrix}$$

$$\text{or: } \theta = \begin{bmatrix} \frac{e^{\mathbf{x}^\top \mathbf{w}_1}}{1 + \sum_{c=1}^{K-1} e^{\mathbf{x}^\top \mathbf{w}_c}} & \dots & \frac{e^{\mathbf{x}^\top \mathbf{w}_{K-1}}}{1 + \sum_{c=1}^{K-1} e^{\mathbf{x}^\top \mathbf{w}_c}} & \frac{1}{1 + \sum_{c=1}^{K-1} e^{\mathbf{x}^\top \mathbf{w}_c}} \end{bmatrix} \begin{matrix} y_1 \\ y_2 \\ y_3 \end{matrix}$$



Connection to neural networks

Multinomial regression:

- Define:

$$\theta = \begin{bmatrix} \frac{e^{\mathbf{x}^\top \mathbf{w}_1}}{1 + \sum_{c=1}^{K-1} e^{\mathbf{x}^\top \mathbf{w}_c}} & \dots & \frac{1}{1 + \sum_{c=1}^{K-1} e^{\mathbf{x}^\top \mathbf{w}_c}} \end{bmatrix}$$

- Cost function is ($z_{i\cdot}$ is one-of- K encoding of y_i)

$$E = - \sum_{i=1}^N \log p(y_i | \mathbf{x}_i) = - \sum_{i=1}^N \sum_{c=1}^K z_{ic} \log \theta_{ic}$$

Multi-class neural network:

- Suppose $\tilde{y}_1, \dots, \tilde{y}_K$ are outputs of a neural network
- Define

$$\theta = \begin{bmatrix} \frac{e^{\tilde{\mathbf{y}}^\top \mathbf{w}_1}}{\sum_{c=1}^K e^{\tilde{\mathbf{y}}^\top \mathbf{w}_c}} & \dots & \frac{e^{\tilde{\mathbf{y}}^\top \mathbf{w}_K}}{\sum_{c=1}^K e^{\tilde{\mathbf{y}}^\top \mathbf{w}_c}} \end{bmatrix}$$

- Cost function is:

$$E = - \sum_{i=1}^N \log p(y_i | \tilde{\mathbf{y}}_i) = - \sum_{i=1}^N \sum_{c=1}^K z_{ic} \log \theta_{ic}$$

Quiz 3, Multinomial Regression (Spring 2016)

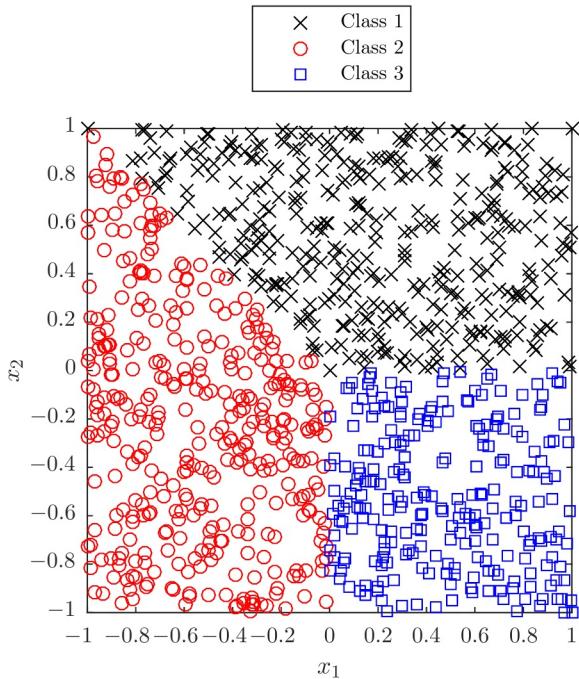


Figure 1: Observations labelled with the most probable class

Consider a multinomial regression classifier for

a three-class problem where for each point $\mathbf{x} = [x_1 \ x_2]^\top$ we compute the class-probability using the softmax function

$$P(\hat{y} = k) = \frac{e^{\mathbf{w}_k^\top \mathbf{x}}}{e^{\mathbf{w}_1^\top \mathbf{x}} + e^{\mathbf{w}_2^\top \mathbf{x}} + e^{\mathbf{w}_3^\top \mathbf{x}}}.$$

A dataset of $N = 1000$ points where each point is labeled according to the maximum class-probability is shown in Figure 1. Which setting of the weights was used?

- A. $w_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$, $w_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, $w_3 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$
- B. $w_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$, $w_2 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$, $w_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$
- C. $w_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, $w_2 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$, $w_3 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$
- D. $w_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$, $w_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, $w_3 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$
- E. Don't know.

Consider for instance the point \mathbf{x} where $x_1 = 0$ and $x_2 = 1$. Then, letting $y_k = \mathbf{w}_k^T \mathbf{x}$, we obtain:

$$A : [y_1 \ y_2 \ y_3] = [-1 \ 1 \ -1]$$

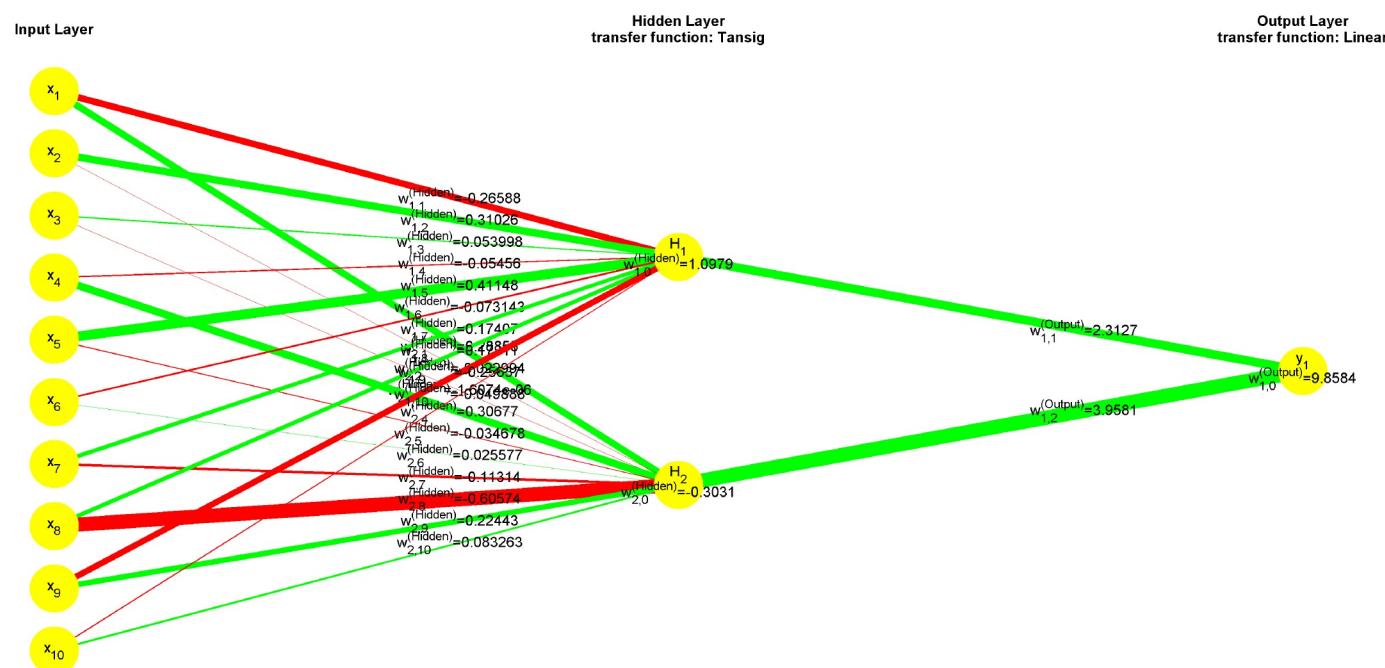
$$B : [y_1 \ y_2 \ y_3] = [-1 \ -1 \ 1]$$

$$C : [y_1 \ y_2 \ y_3] = [1 \ -1 \ -1]$$

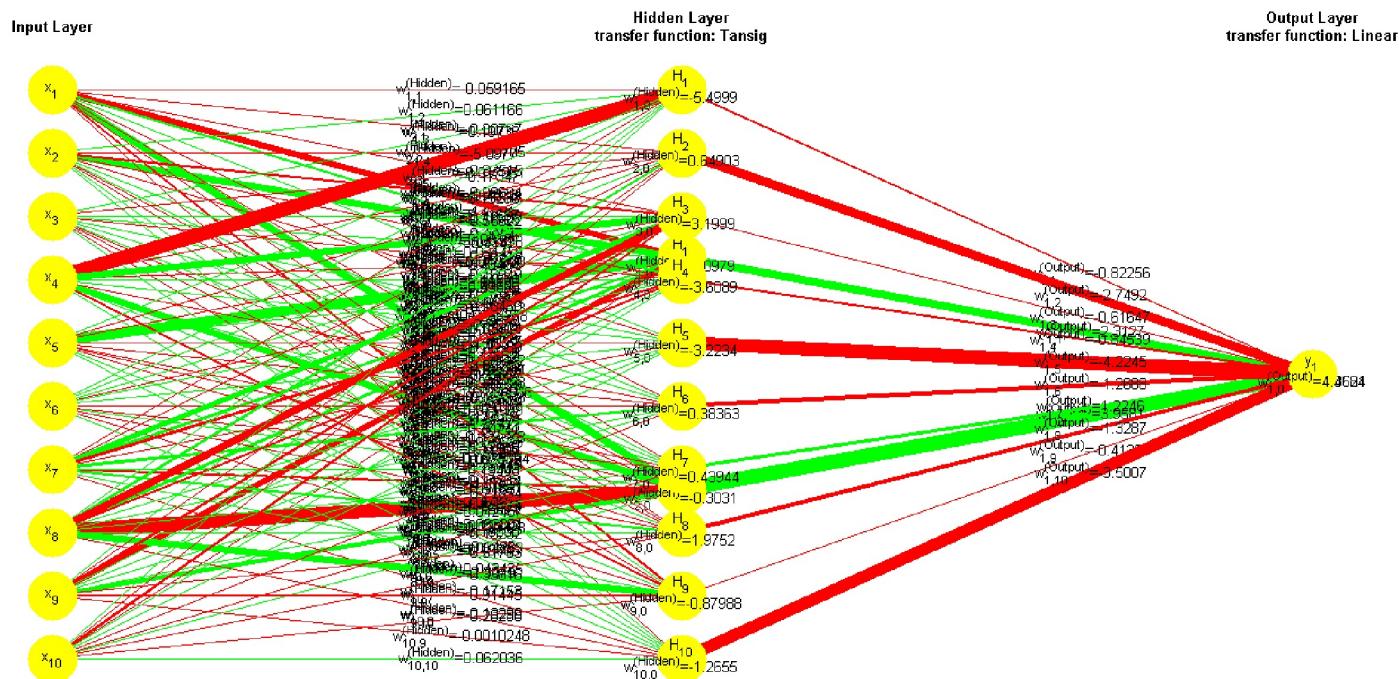
$$D : [y_1 \ y_2 \ y_3] = [-1 \ 1 \ 1]$$

Next, since the multinomial regression function preserves order we need only consider the maximal value. Accordingly the point \mathbf{x} is only classified to the correct class 1 for option C .

Interpreting neural networks can be difficult



Interpreting neural networks can be difficult



Resources

<https://www.youtube.com> Excellent video resource explaining the concepts behind neural networks

(https://www.youtube.com/watch?v=aircAruvnKk&list=PLZHQB0b0WTQDNU6R1_67000Dx_ZCJB-3pi)

<http://playground.tensorflow.org> Sleek interactive neural network example where you can examine the effect of different number of hidden neurons, activation functions, and many other things on training (<http://playground.tensorflow.org/>)

<https://www.tensorflow.org> Most popular and well-documented deep learning framework. While well documented, notice it requires some python knowledge (<https://www.tensorflow.org/>)

<https://pytorch.org> Upcoming (and in some ways slightly simpler) framework for deep learning; alternative to tensorflow

(<https://pytorch.org/>)