



COMP0197 Applied Deep Learning

Andre Altmann

Department of Medical Physics and Biomedical Engineering
Centre for Medical Image Computing (CMIC)

a.altmann@ucl.ac.uk

London, January 11th 2024

My Background



Master's in Computer Science

- Speech and Language technologies
- Machine learning

PhD in Computational Biology

- Machine Learning ... applied to biological problems

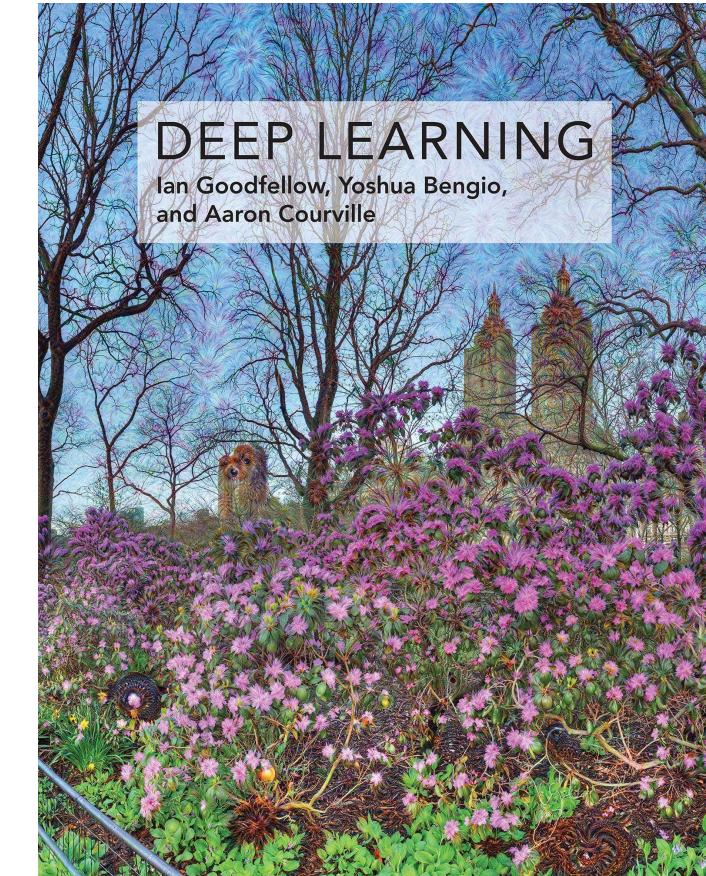
Associate Professor in Medical Physics Department

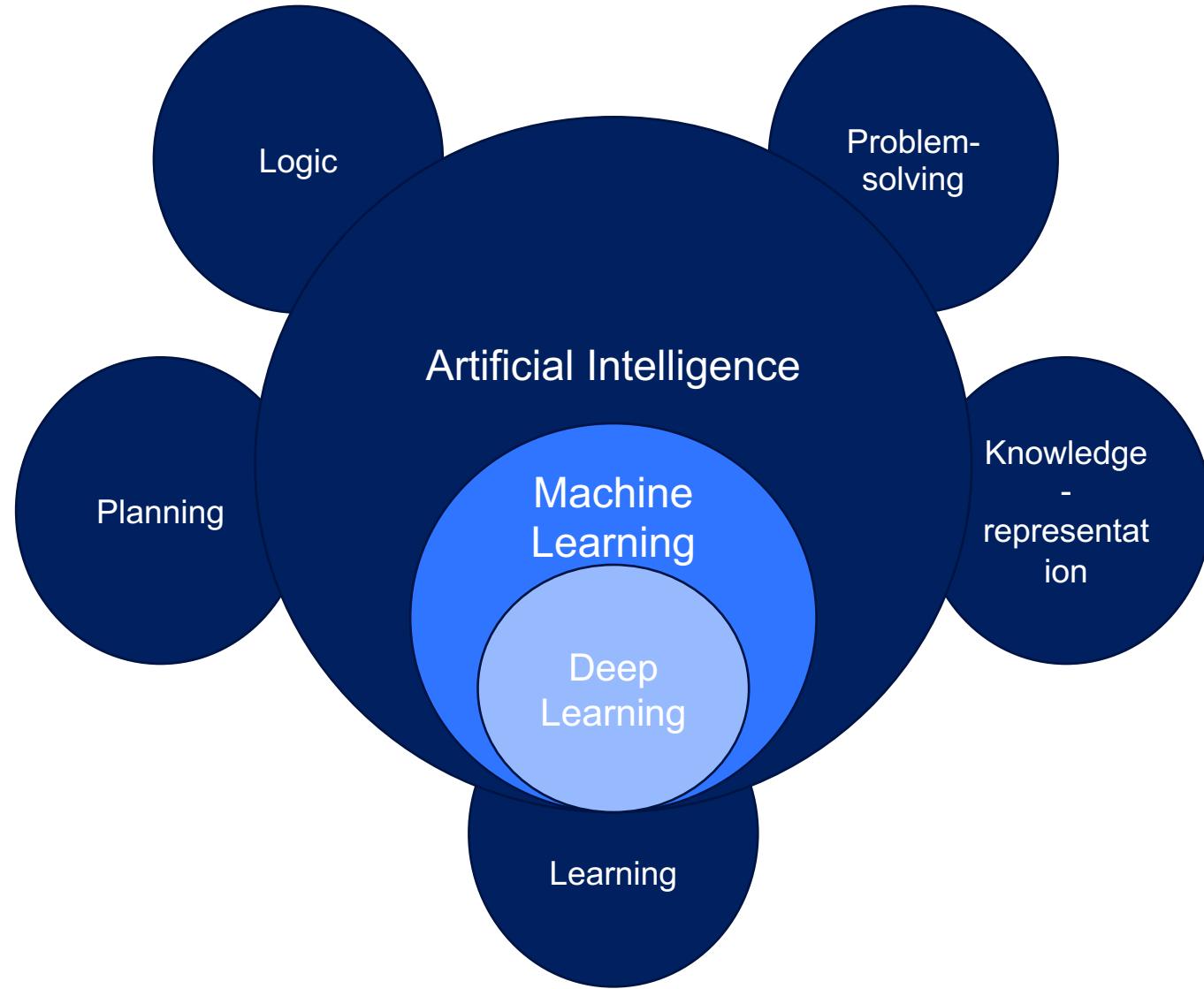
- Linking medical imaging data to genetics etc. ...
- ... using Machine Learning

A refresher on linear algebra and Probability Theory can be found in Chapters 2 and 3

Today

- T, E, P
- Linear regression
- Overfitting/underfitting/capacity
- Hyperparameters
- Bias-Variance
- Making decisions





Task
Experience
Performance

A definition of (machine) learning

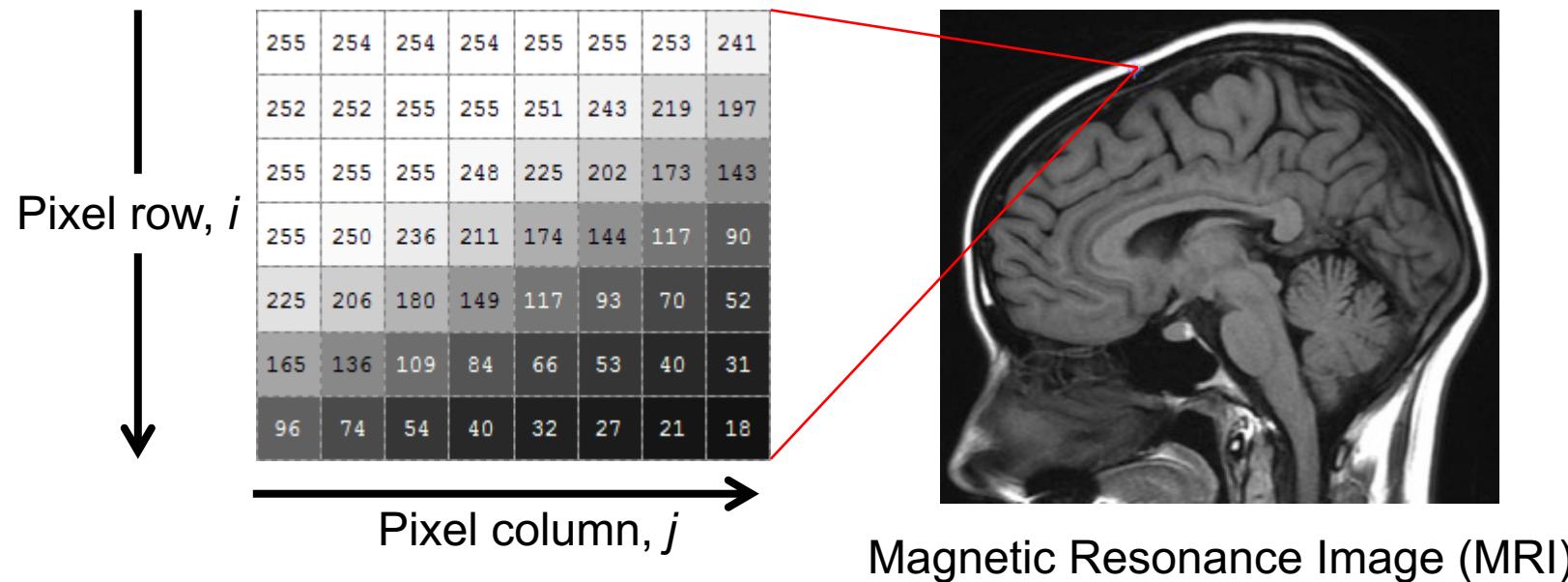


“A computer program is said to learn from experience E with respect to some class of task T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .”

- Mitchell (1997)

Some definitions

- E (experience) comes in the form of **examples**
- Examples are collections of **features**



- Formally: $x \in \mathbb{R}^n$
 - Each entry of the vector x_i is a **feature**

Classification

- Specify to which of k categories an input belongs to
 $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$

Many tasks are ‘binary’ problems (yes vs. no; $k=2$)

- One-hot encoding: $(0,0,1,0) \leftrightarrow 3$

Regression

- Predict a numerical value from the input:

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Further groups of tasks



Transcription: produce text describing a picture, transcribe spoken words to written text, ...

Machine translation: translate text from one language to another

Structured output: Output is a vector rather than a scalar

Density estimation: Learn $p_{model}: \mathbb{R}^n \rightarrow \mathbb{R}$, where p_{model} is viewed as a probability density function

...

The performance measure P



- The performance measure P is specific to the task T
- Often used:
 - **Accuracy:** “percentage correct”
 - **Error rate:** “percentage incorrect” (0-1-Loss)
 - ...
- P can sometimes be hard to ‘quantify’:
 - Translation: need to consider word ordering etc.
 - Density estimation

The performance measure P



- The performance measure can massively influence what the algorithm learns!
- A regression problem, typically we optimize the “mean squared error”

$$\text{MSE} = \frac{1}{m} \sum_i (\hat{y}_i - y_i)^2$$



- MSE is sensitive to ‘outliers’
- Better mean absolute error: $\text{MAE} = \frac{1}{m} \sum_i |\hat{y}_i - y_i|$

Interested in how well the algorithm performs on new, unseen data:

generalization

= good performance P on an independent **test set** of data, i.e., not seen during training time

The experience E



Typically, the *Experience* is a **dataset** comprising many examples aka **data points**

Two board categories

Unsupervised

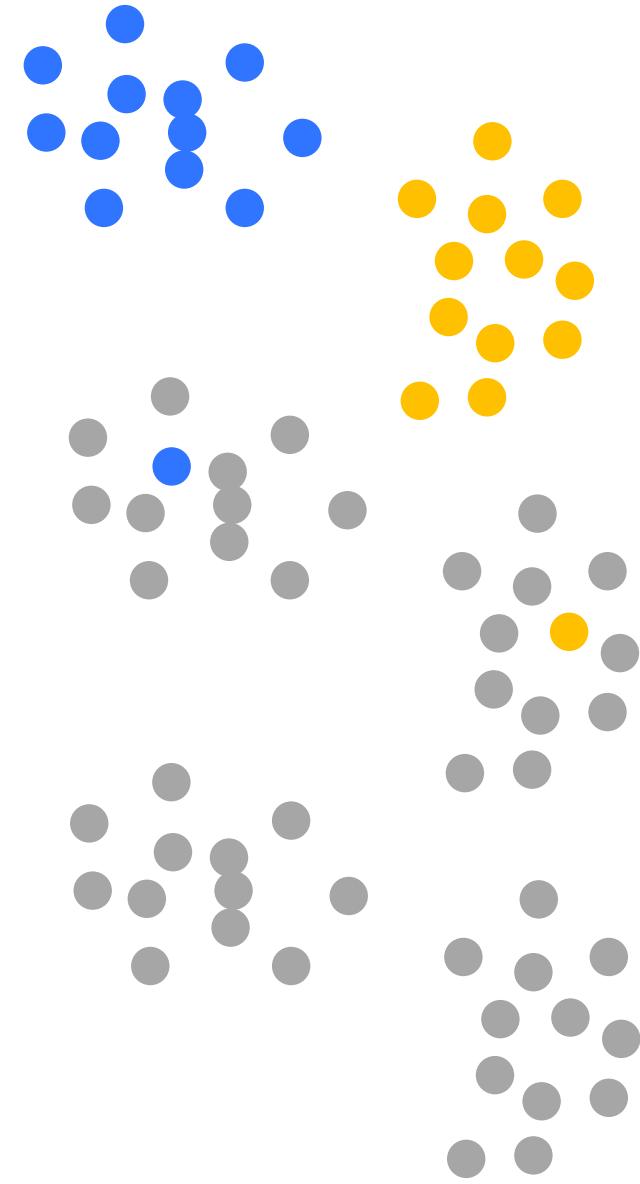
- contains many **features** and the task is to extract interesting properties

Supervised

- contains **features** as well as an **output label/target**

Broad categories of ML

- Supervised Learning
 - We are provided a label/target value
 - E.g. support vector machines
- Semi-supervised learning
 - Exploit distributions etc. of unlabeled data to improve supervised learning
- Unsupervised Learning
 - No labels used, aim to find a pattern
 - E.g., clustering



The experience E



- Often our examples are vectors of the same length
- **Design matrix** $X \in \mathbb{R}^{m \times n}$, m : samples, n : features
- More general: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
 - Each x can have a different dimension (e.g., images, sentences)

Linear regression

Linear regression

- Input $x \in \mathbb{R}^n$
- Output $y \in \mathbb{R}$
- The model assumes that y is a linear combination of the values in x :
$$\hat{y} = w^T x.$$
- $w \in \mathbb{R}^n$: model parameters learnt from the training data

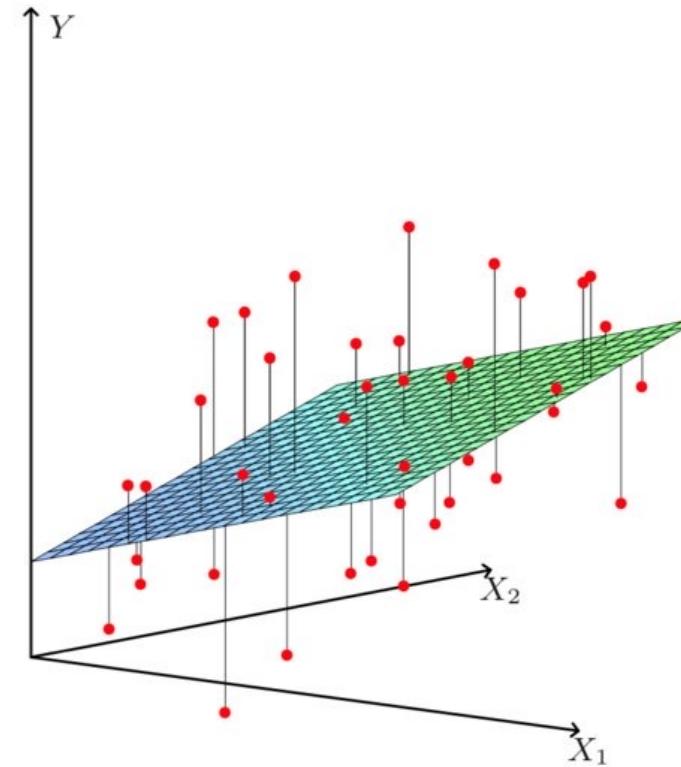


FIGURE 3.1. Linear least squares fitting with $X \in \mathbb{R}^2$. We seek the linear function of X that minimizes the sum of squared residuals from Y .

Linear regression

- T : predict y from x by computing
 $\hat{y} = \mathbf{w}^T \mathbf{x}$.
- P : “*distance from observation y to prediction \hat{y} .*”
 - Assume we have a test set of m samples $X^{(\text{test})}$ along with the targets $y^{(\text{test})}$.
 - Mean squared error:

$$\begin{aligned}\text{MSE}_{\text{test}} &= \frac{1}{m} \sum_i (\hat{y}^{(\text{test})} - y^{(\text{test})})_i^2 \\ &= \frac{1}{m} \|\hat{y}^{(\text{test})} - y^{(\text{test})}\|_2^2\end{aligned}$$

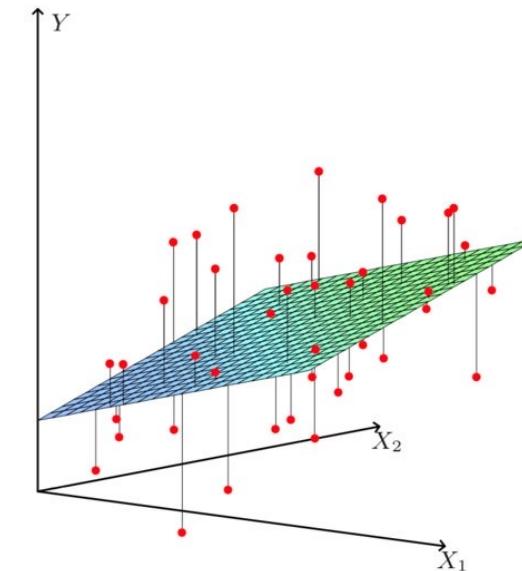


FIGURE 3.1. Linear least squares fitting with $X \in \mathbb{R}^2$. We seek the linear function of X that minimizes the sum of squared residuals from Y .

Find \mathbf{w}^T that minimize MSE_{test}

Linear regression



During training we have the training data $\{\mathbf{X}^{(\text{train})}, \mathbf{y}^{(\text{train})}\}$.

We devise an ML algorithm by minimizing $\text{MSE}_{\text{train}}$

- Has a '*closed form*' solution:

$$\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = 0.$$

('set the derivative of $\text{MSE}_{\text{train}}$ to 0')

Linear regression



$$\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = 0$$

$$\Rightarrow \nabla_{\mathbf{w}} \frac{1}{m} \left\| \hat{\mathbf{y}}^{(\text{train})} - \mathbf{y}^{(\text{train})} \right\|_2^2 = 0$$

$$\Rightarrow \nabla_{\mathbf{w}} \frac{1}{m} \left\| \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right\|_2^2 = 0$$

$$\Rightarrow \nabla_{\mathbf{w}} (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})})^T (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}) = 0$$

$$\Rightarrow \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^{(\text{train})T} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^{(\text{train})T} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})T} \mathbf{y}^{(\text{train})}) = 0$$

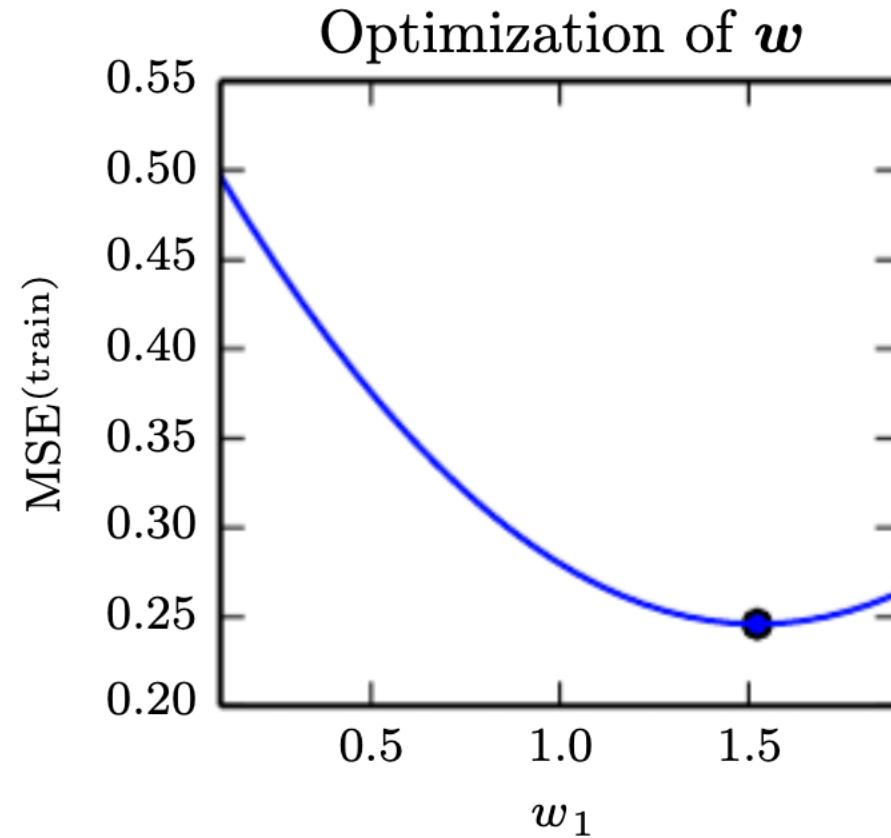
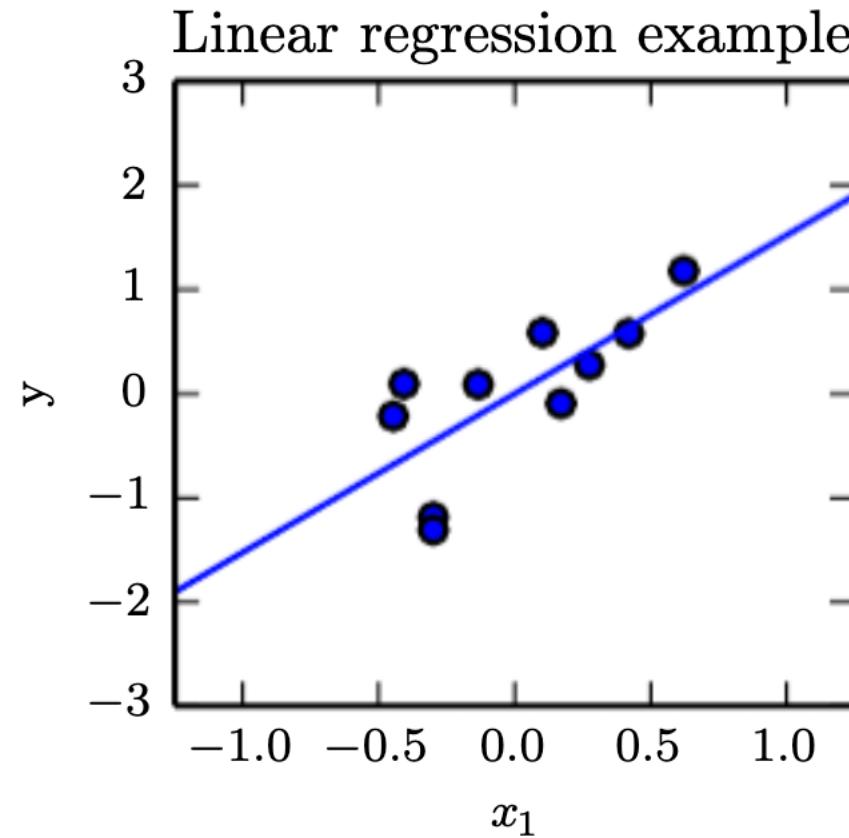
$$\Rightarrow 2\mathbf{X}^{(\text{train})T} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{X}^{(\text{train})T} \mathbf{y}^{(\text{train})} = 0$$

$$\Rightarrow \mathbf{w} = (\mathbf{X}^{(\text{train})T} \mathbf{X}^{(\text{train})})^{-1} \mathbf{X}^{(\text{train})T} \mathbf{y}^{(\text{train})} \quad \text{Normal equation}$$

Linear regression:

$$\hat{\mathbf{y}} = \mathbf{w}^T \mathbf{x} + b.$$

Linear regression



Essence of (supervised) machine learning

- We have a model that represents types of functions and is parameterized by θ

$$f(\mathbf{x}; \boldsymbol{\theta})$$

- We find the setting θ^* that minimizes the **Loss** (makes the fewest errors) on our training data:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y_i)$$

Overfitting, Underfitting & Capacity

Overfitting/Underfitting



- In practice we optimize the **training error**
- However, we aim to minimize the **test error** or **generalization error**
- Measured on a test set $\{X^{(\text{test})}, y^{(\text{test})}\}$:

$$\frac{1}{m^{\text{test}}} \|X^{(\text{test})} \mathbf{w} - y^{(\text{test})}\|_2^2$$

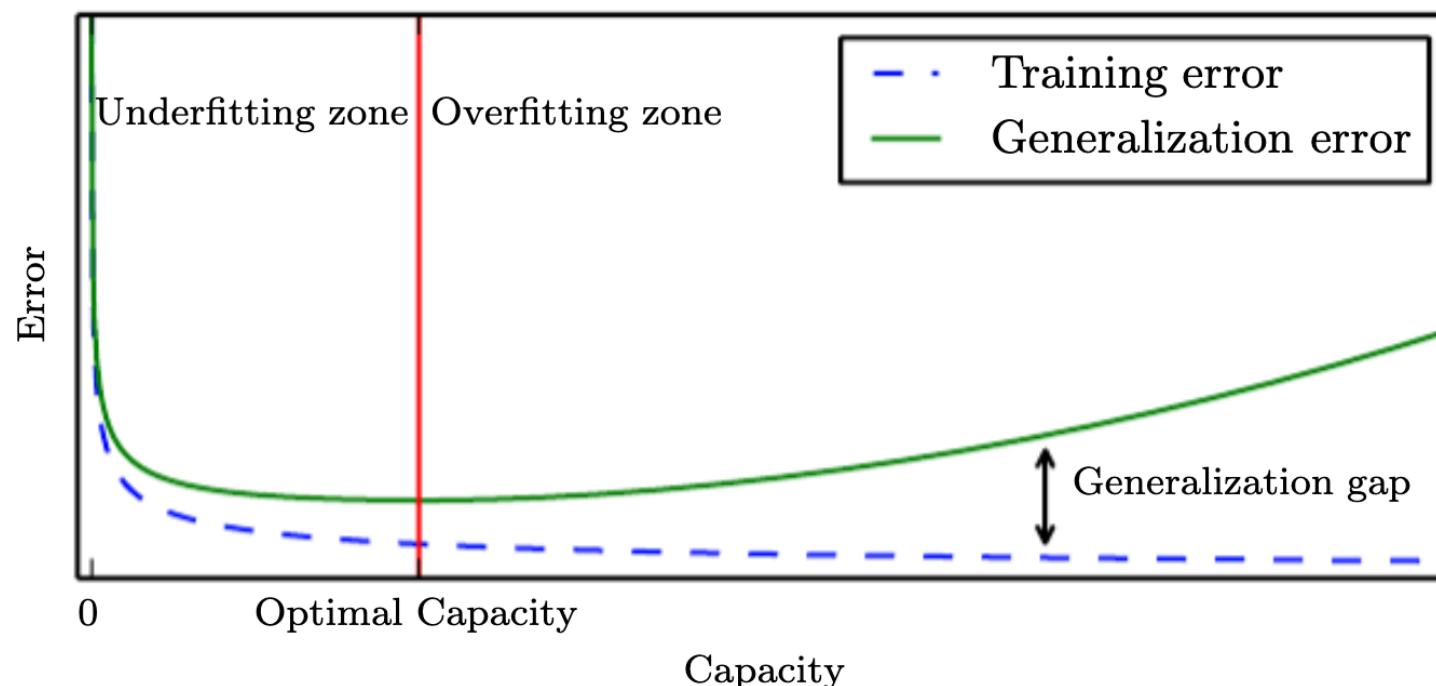
- Works: if training and test data were generated by the same data-generating process: $p(x, y)$.
- **Realistic assumption?**

- Test sets need to be largish
- Accuracy $\frac{1}{m} \sum_i^m \delta_{y_i=\hat{y}_i}$ is essentially a ‘mean’
- Standard error of the mean:
$$SE = \sigma / \sqrt{m}$$
- Thus, accuracy estimate (with 95% CI):
$$[Acc - 1.96 * SE, Acc + 1.96 * SE]$$

Thus, if your test set is small, you have high uncertainty on your Acc (or other performance measure)

Overfitting/Underfitting

- Gap between training and test error
 - Overfitting: model too tuned to training data and performance on test data declines
 - Underfitting: model does not use full potential of the training data



- Model “*flexibility*” as in ability to fit a wide variety of functions

- Linear model rather ‘inflexible’:

$$\hat{y} = b + wx$$

- Adding ‘powers’ adds flexibility:

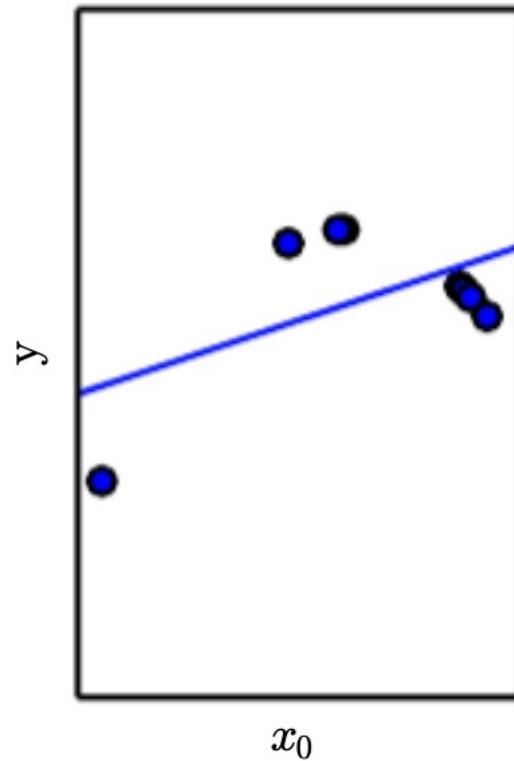
$$\hat{y} = b + w_1x + w_2x^2$$

$$\hat{y} = b + \sum_i w_i x^i$$

- The power i becomes a model parameter

Capacity

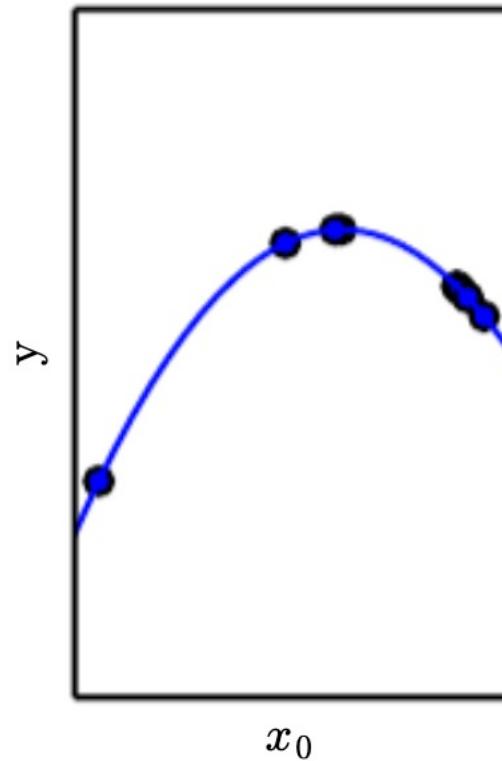
Underfitting



x_0

$i=1$

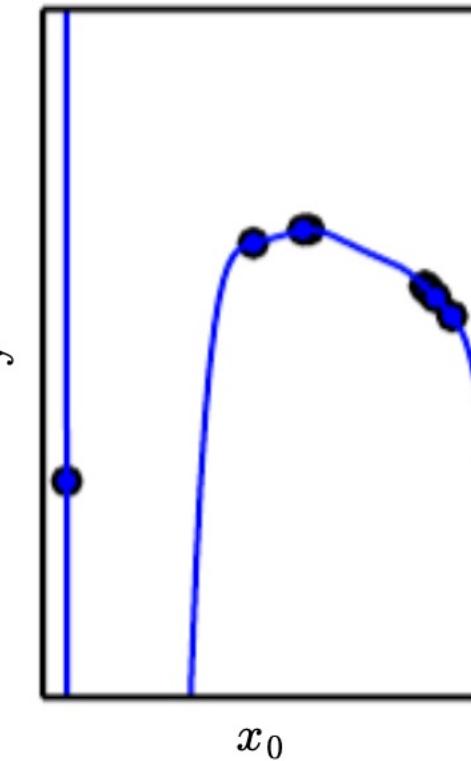
Appropriate capacity



x_0

$i=2$

Overfitting



x_0

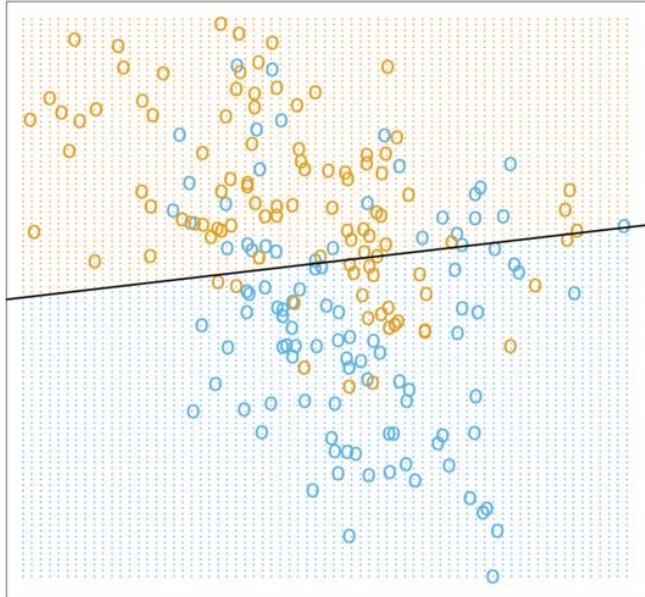
$i=9$

Representational capacity

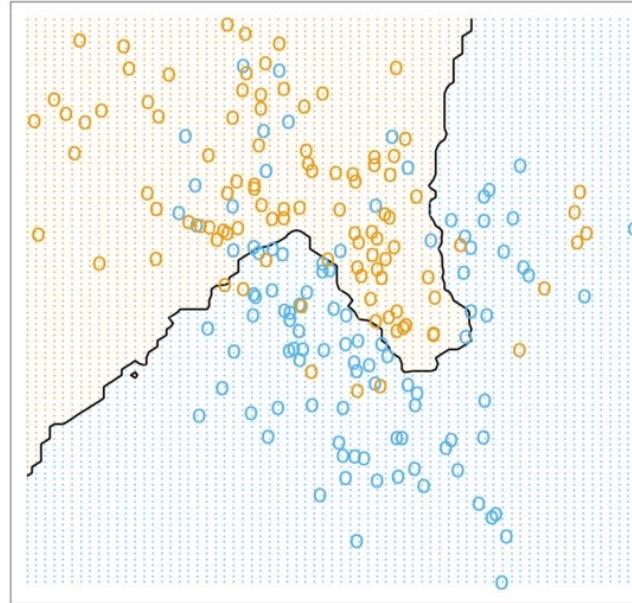
Occam's razor

Capacity

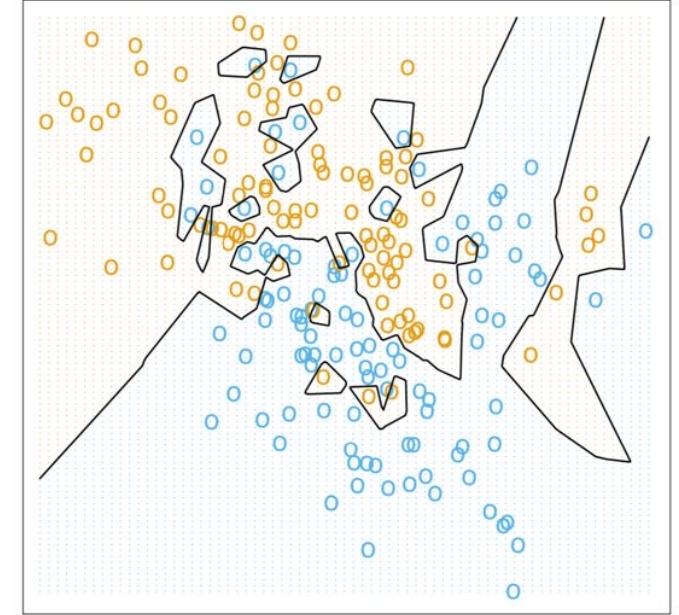
Linear Regression of 0/1 Response



15-Nearest Neighbor Classifier



1-Nearest Neighbor Classifier



Nearest Neighbor Classifier

$$\hat{y} = y_i$$
$$i = \operatorname{argmin} \|X_i - x\|_2^2$$

Regularization & Hyperparameter tuning

No Free Lunch Theorem



“Averaged over all possible data-generating distributions, every classification algorithm has the same error rate when classifying previously unobserved points.”

- Wolpert (1996)

Regularization



- Using background knowledge in algorithm design improves the model
 - E.g., we know the data follows a quadratic function etc.
 - Dependencies in real life complex
 - In practice we constrain the possible solutions

- Previously we optimized only P :

$$J(\mathbf{w}) = \text{MSE}_{\text{train}}$$

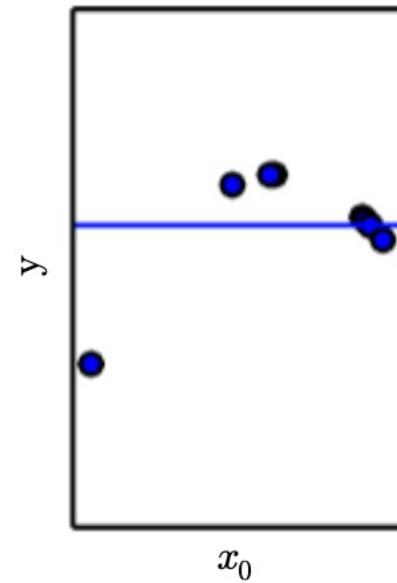
- Now we ‘punish’ model complexity:

$$J(\mathbf{w}; \lambda) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$$

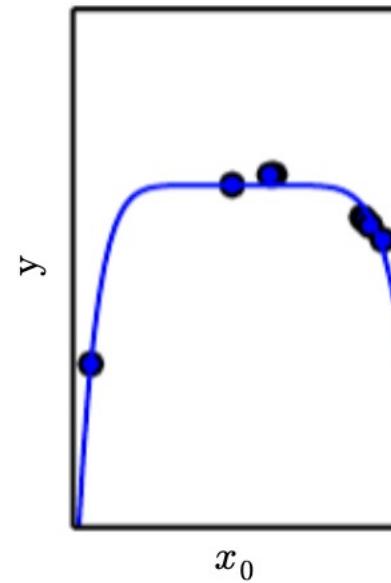
*L² –norm
weight decay*

Regularization

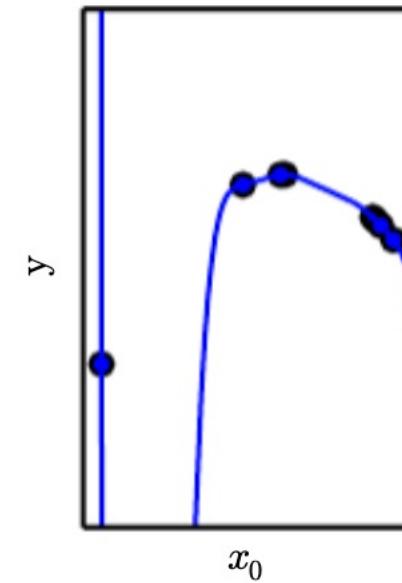
Underfitting
(Excessive λ)



Appropriate weight decay
(Medium λ)



Overfitting
($\lambda \rightarrow 0$)

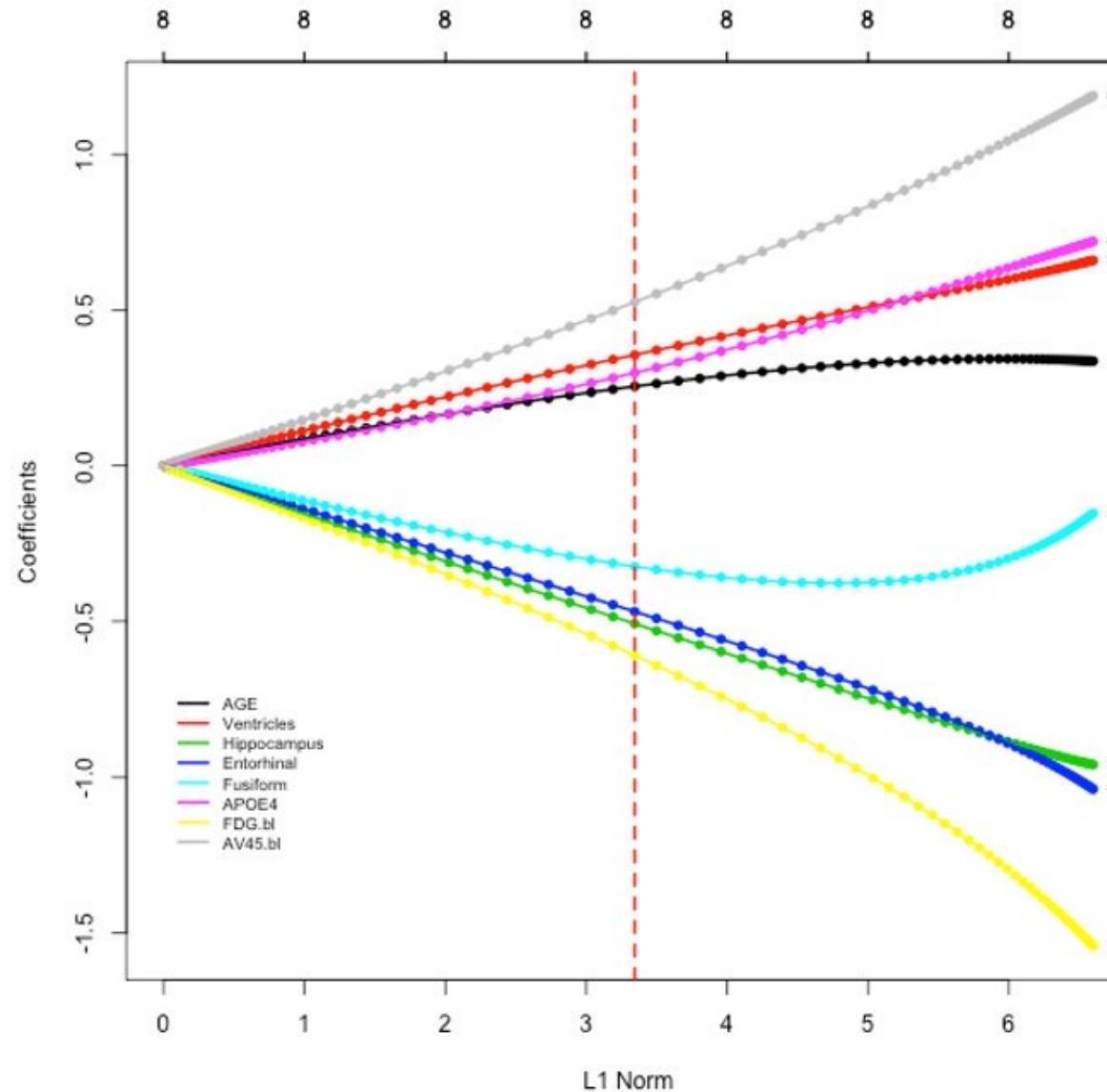


$$\hat{y} = b + \sum_i^9 w_i x^i$$

$$\Omega(\mathbf{w}) = \mathbf{w}^T \mathbf{w}$$

Regularizer

Regularization – Ridge regression



Tuning hyperparameters



- *Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.*
- In the example before λ is a parameter that controls **capacity: hyperparameter**
- How to tune the hyperparameter?

Tuning hyperparameter

- An ‘independent’ dataset: **validation set**
- “training set for hyper parameters”



80%

20%

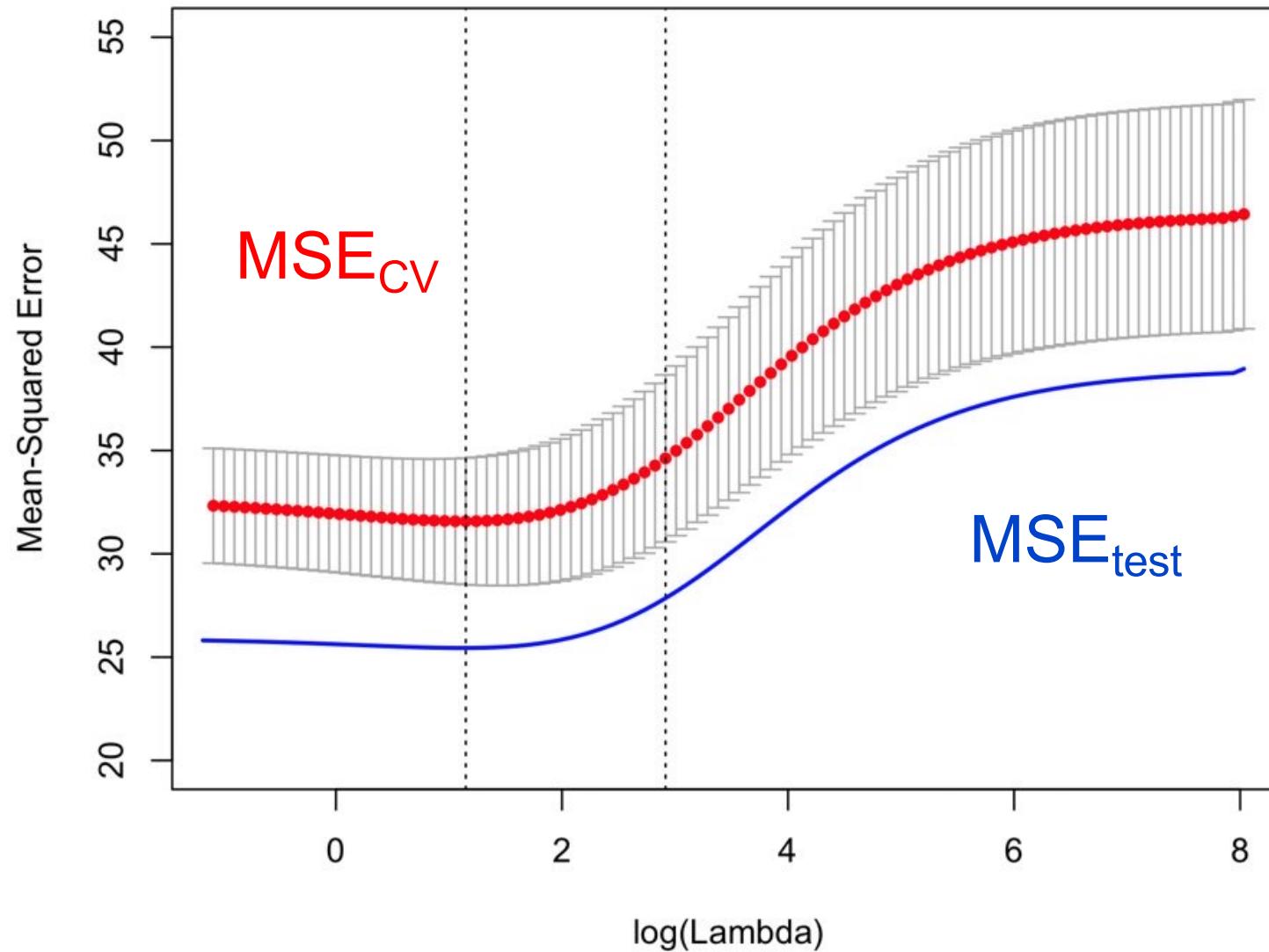
Cross-validation (CV)

- Train-validation split works well when we are ‘data rich’
- For small data set: split repeatedly
- Split into k non-overlapping sets (‘ k -fold’)
 - Rotate role of validation set



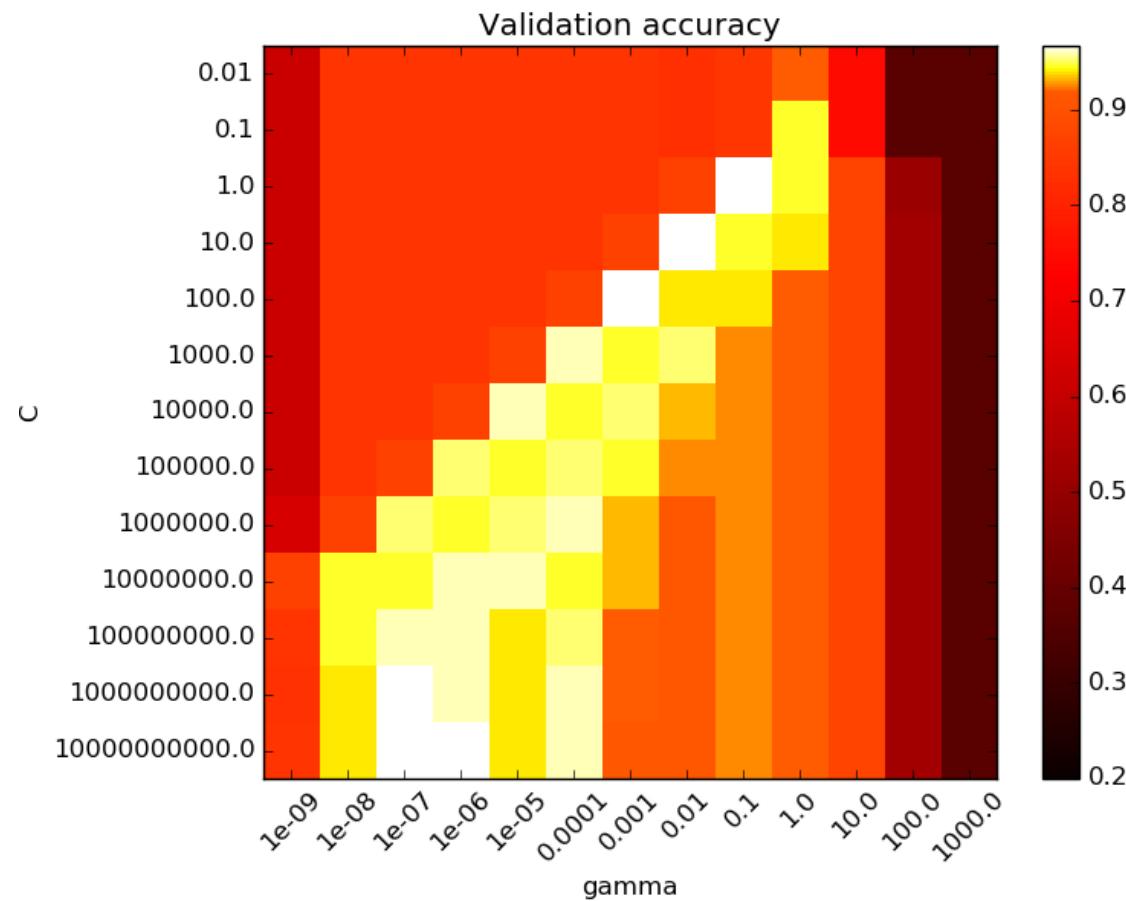
5-fold cross-validation

Cross-validation (CV)



Grid search

- More than one hyperparameter?



Bias Variance Trade-off

Bias-variance-tradeoff



- One can decompose MSE into Bias and variance

$$\begin{aligned}\text{MSE} &= \mathbb{E}[(\hat{\theta}_m - \theta)^2] \\ &= \text{Bias}(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m)\end{aligned}$$

- Bias: “Error made by assumption.”
- Variance: “Error made b/c adopting to data.”

Point Estimates



- True value: θ
- Estimate from data: $\hat{\theta}$
- Let $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ be a set of m i.i.d. points
- Then: $\hat{\theta}_m = g(x^{(1)}, \dots, x^{(m)})$
- $\hat{\theta}_m$ is a point estimate of θ , which generated the data
- Example: $\mathcal{N}(x^{(i)}; \mu, \sigma^2)$ generates data; $\theta = (\mu, \sigma)$
by analyzing m datapoints we get $\hat{\theta}_m = (\hat{\mu}_m, \hat{\sigma}_m)$

- Estimates deviate from the ‘ground truth’: bias
 $\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta$
- **Unbiased:** if $\text{bias}(\hat{\theta}_m) = 0$
- **Asymptotically unbiased:** $\lim_{m \rightarrow \infty} \text{bias}(\hat{\theta}_m) = 0$

Example: Gaussian distrib.



$$p(x^{(i)}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x^{(i)} - \mu)^2}{\sigma^2}\right)$$

- Sample mean: $\hat{\mu}_m = \frac{1}{m} \sum_i x^{(i)}$
- Bias: $\text{bias}(\hat{\mu}_m) = \mathbb{E}[\hat{\mu}_m] - \mu = \mathbb{E}\left[\frac{1}{m} \sum_i x^{(i)}\right] - \mu$
$$= \left(\frac{1}{m} \sum_i \mathbb{E}[x^{(i)}]\right) - \mu$$
$$= \left(\frac{1}{m} \sum_i \underline{\mu}\right) - \mu = \underline{\mu} - \underline{\mu} = 0$$

Example: Gaussian distrib.



$$p(x^{(i)}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x^{(i)} - \mu)^2}{\sigma^2}\right)$$

- Sample variance: $\hat{\sigma}^2_m = \frac{1}{m} \sum_i (x^{(i)} - \hat{\mu}_m)^2$
- Bias: $\text{bias}(\hat{\sigma}^2_m) = \mathbb{E}[\hat{\sigma}^2_m] - \sigma^2$
- $\mathbb{E}[\hat{\sigma}^2_m] = \mathbb{E}\left[\frac{1}{m} \sum_i (x^{(i)} - \hat{\mu}_m)^2\right] = \dots = \frac{m-1}{m} \sigma^2$

Sample variance is biased. That is why we use $\frac{1}{m-1} \sum_i (x^{(i)} - \hat{\mu}_m)^2$ instead.

Variance of an Estimator



- “How much does the estimator vary as a function of the data sample?”

$$\text{Var}(\hat{\theta})$$

- E.g., an estimator $g(x^{(1)}, x^{(2)}, \dots, x^{(m)}) = c$ returning a constant has variance 0

Bias-variance-tradeoff



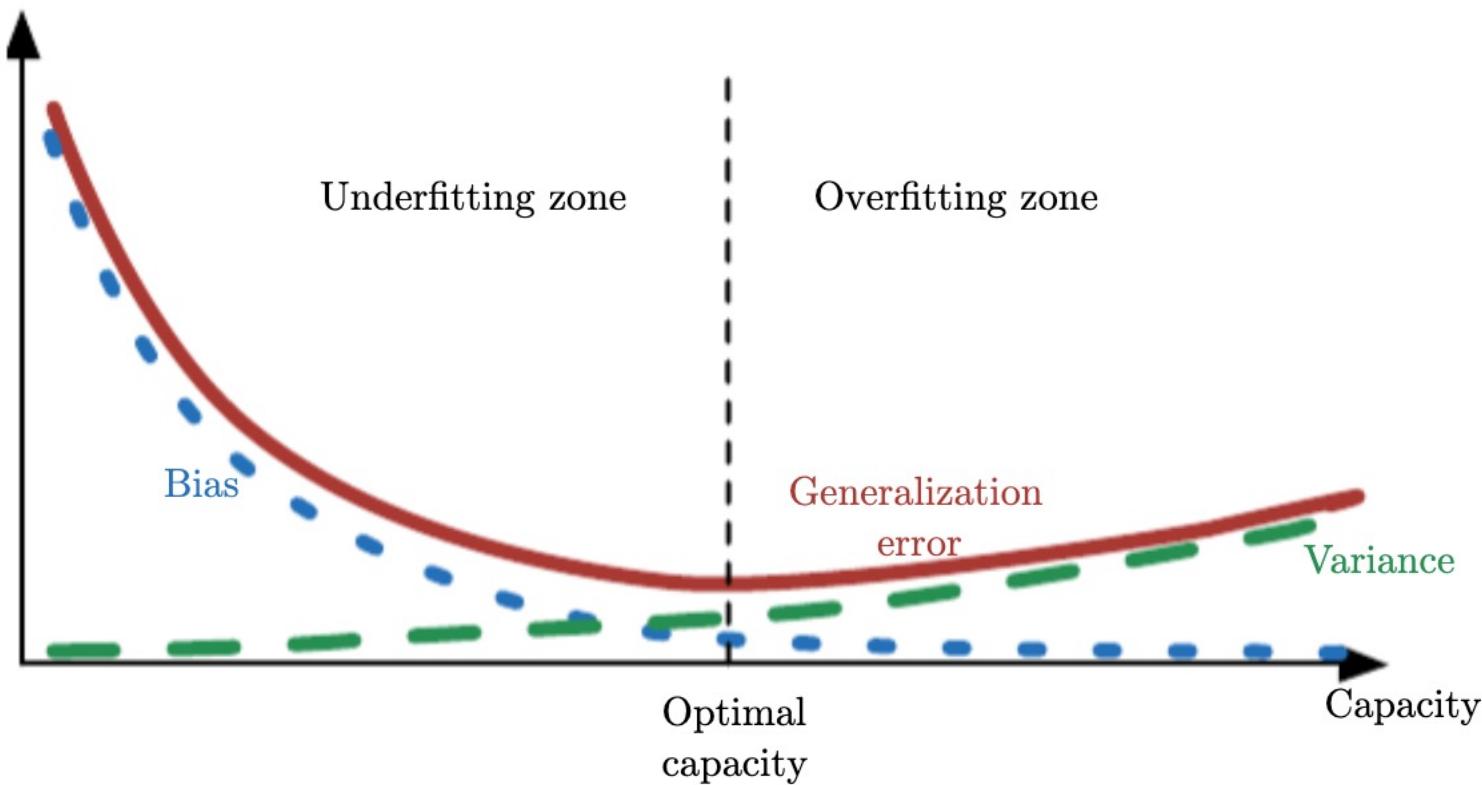
- One can decompose MSE into Bias and variance

$$\begin{aligned}\text{MSE} &= \mathbb{E} [(\hat{\theta}_m - \theta)^2] \\ &= \text{Bias}(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m)\end{aligned}$$

- Bias: “Error made by assumption.”
- Variance: “Error made b/c adopting to data.”

$$J(\mathbf{w}) = \frac{1}{m} \sum_i (\mathbf{x}^{(\text{train})} \mathbf{w} - y^{(\text{train})})_i^2 + \lambda \mathbf{w}^T \mathbf{w}$$

Bias-variance-tradeoff



$$J(\mathbf{w}) = \frac{1}{m} \sum_i (\mathbf{x}^{(\text{train})} \mathbf{w} - y^{(\text{train})})_i^2 + \lambda \mathbf{w}^T \mathbf{w}$$

Making decisions

Discriminant functions

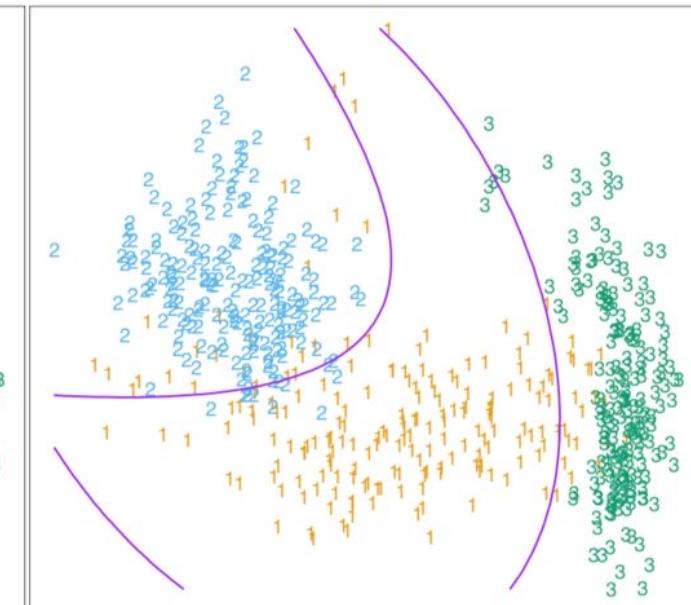
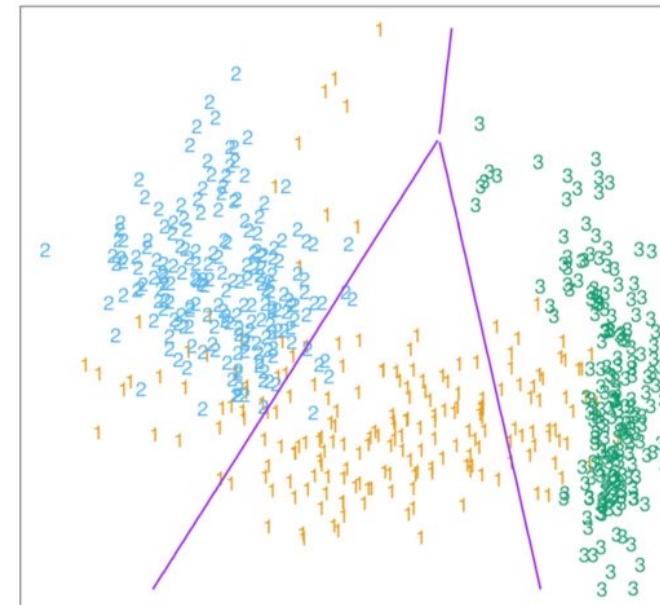
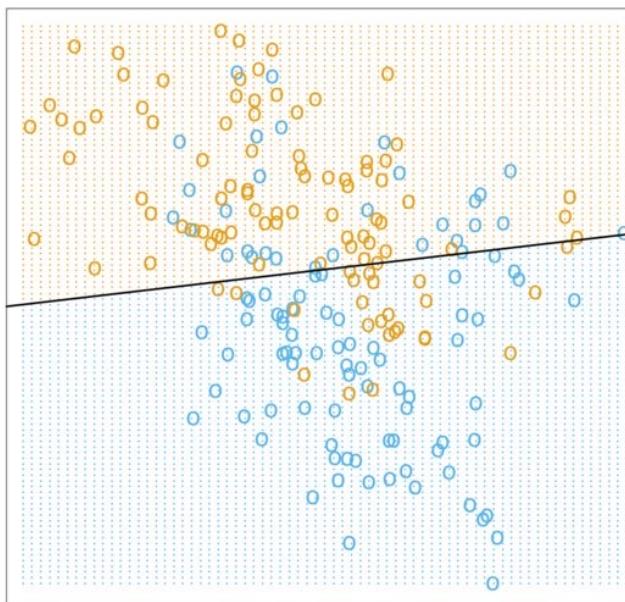


- In **Regression** we have functions $f: \mathbb{R}^n \rightarrow \mathbb{R}$
- In **Classification** our target is $\in \{1, \dots, k\}$
- Methods typically model continuous *discriminant function* for each class: $\delta_k(x)$
- Predicted class is the one with the largest value for $\delta_k(x)$
- E.g., $P(y = k|x)$ is a discriminant function

Decision boundaries

- Where the support for class k and k' is the same
- $\delta_k(x) = \delta_{k',\ell}(x)$

Linear Regression of 0/1 Response



Bayes Theorem

- Very important and useful:

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

- Applied to classification:

$$P(y = k|x) = \frac{P(x|y = k) * P(y = k)}{P(x)}$$

likelihood

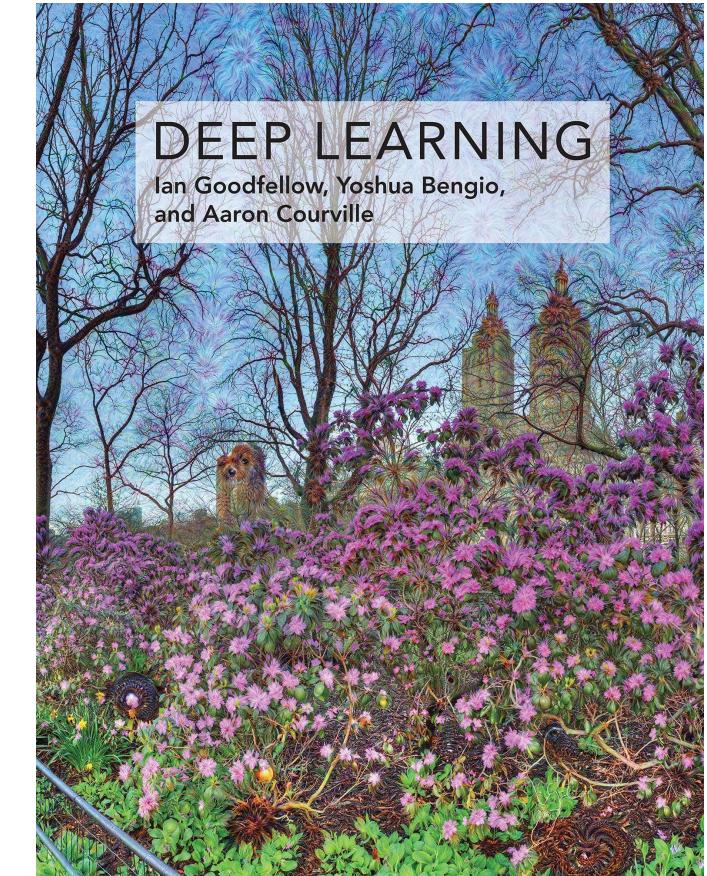
Class prior

support

A refresher on linear algebra and Probability Theory can be found in Chapters 2 and 3

Today

- T, E, P
- Linear regression
- Overfitting/underfitting/capacity
- Hyperparameters
- Bias-Variance
- Making decisions



Thank you!