

ML4P

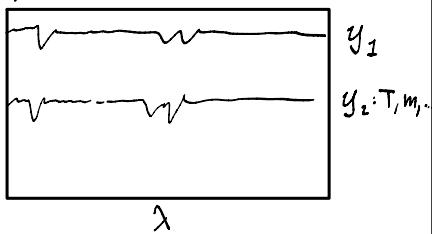
Lecture 2

Types of ML:

- Nearest neighbor
- Linear regression



Spectra



$\dim(y_i) = 7 \times 3$ pixels per spectrum

$n = \#$ spectra in training:
(x_i)

Introduction of symbols

Supervised (Training):

Training set: $\begin{cases} \text{Features: } x_i \\ \text{Labels: } y_i \end{cases}$ bvr 792

Train $f(x_i, \hat{W}) \sim y_i$,

Test data: x_i^*, y_i^* & x_i, y_i to find
the weights, \hat{W} .

$$\mathbf{X} = [x_1, \dots, x_n]^T$$

$y_i = \hat{a} \cdot \hat{x}_i + \text{noise}$

$\hat{a} = (\mathbf{X}^T C^{-1} \mathbf{X})^{-1} \mathbf{X}^T C^{-1} \mathbf{Y}, \quad \mathbf{Y} = [y_1, y_2, \dots, y_n]^T$

Fitting, in closed form:
Note: Not lin. regression

► Fitting: I want to know a
Regression: We want to predict
 y^* .

- seek \hat{a} s.t. $E[\|\hat{a} \cdot \vec{x}^* - y^*\|]$ is small.

- features are independent

- test data are drawn from the same distribution
as training

- is a linear relation

↳ Between \vec{a} , \vec{x} ?

Then: $\hat{a} \leftarrow \underbrace{(X^T X)^{-1} X^T Y}_{\substack{p \times p \\ p \times n \\ n \times 1}}$ delivers the MVUE

↳ symmetric, real, non-negative eigenvalues.

When: $p \rightarrow \infty$: $\hat{a} \leftarrow X^T (X X^T)^{-1} Y$

Uml python

`linalg()`, `@`, `solve()`

never use: `inv`.

Lecture 3

Linear Regression Revisited

(2nd) ML definition: We care about data, not parameters prediction at

Prediction

In fact, we will choose a larger number of parameters than data.

Important Lin. Alg.

- Rank
- SVD
- Condition #
- Non-negative definite.

$$\Leftrightarrow \lambda \geq 0$$

Recall $\hat{a} = \text{vector of parameters}$

$$= (\underbrace{X^T X}_{p \times n \quad n \times p})^{-1} X Y$$

$\underbrace{}$

$$\text{rank } \leq \min(h, p)$$

Since we only care about the prediction of data:

$$\left\{ \begin{array}{l} \hat{y}^* = X^* \hat{a} \\ \hat{y}^* = X^* (X^T X)^{-1} X Y \quad n > p \\ \hat{y}^* = X^* X^T (X X^T)^{-1} Y \quad n < p \end{array} \right\}$$

Pseudo inverse: $Y^* = X^* X^+ Y$

X^+ definition: $\forall Y = X^+ Y = \underset{\hat{a}}{\arg \min} \| Y - X \cdot \hat{a} \|_2^2 + \lambda \| \hat{a} \|_2^2$

if $n > p$ this cannot be solved

Regularisation parameter

if $p > n$ no real λ

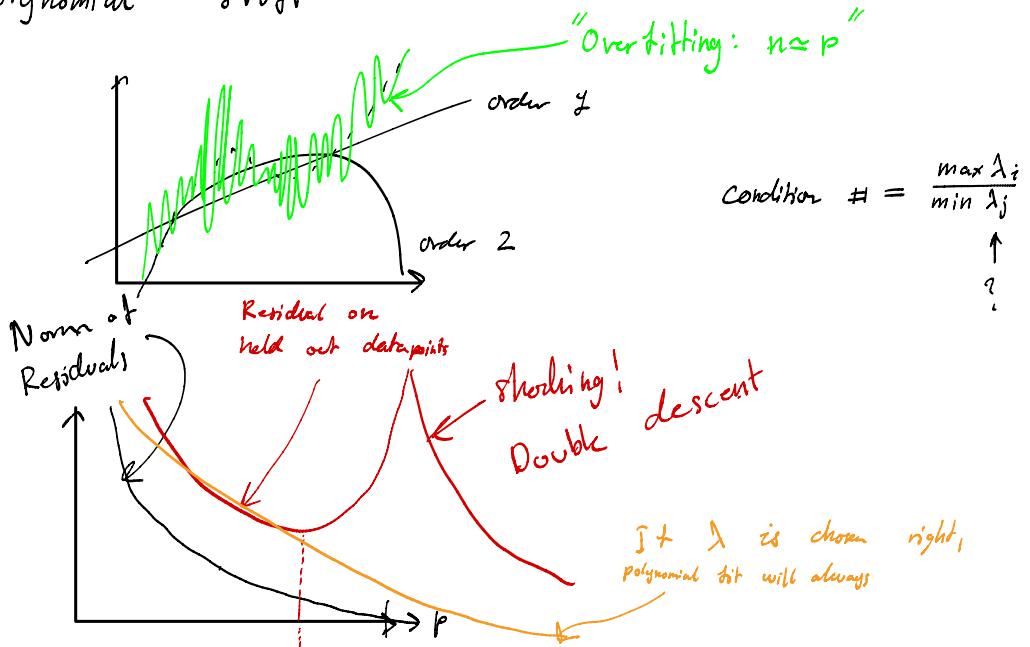
$$\|q\|_a^b \equiv \left[\sum_i |q_i|^a \right]^{b/a}$$

Norm since a can in general be a matrix/tensor.

$$\Lambda = \lambda I_{p \times p}, \quad \hat{y}^* = \lim_{\lambda \rightarrow 0} X^* \underbrace{(X^T X + \Lambda)}_{p \times p}^{-1} X^T Y$$

$$\hat{y}^* = \lim_{\lambda \rightarrow 0} X^* \Lambda^{-1} X^T \underbrace{(X \Lambda^{-1} X^T + I)}_{n \times n}^{-1} Y$$

Polynomial Stuff



Lecture 4

So far only lin. reg., how do we evaluate the quality of a model?

- Train, validation, test, cross validation

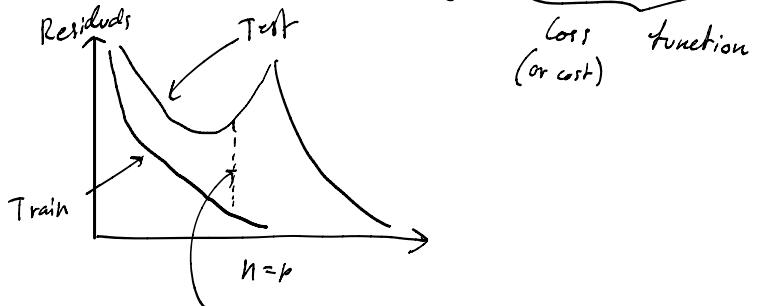
- Feature importance

Supervised: ML is about prediction.

A good (ML) model predicts held out data well.

Last time (No distribution): $\hat{a} = \underset{a}{\operatorname{argmin}} \left[\|Y - X \cdot a\| + \lambda \|a\| \right]$

Overttting:



Difference is large? Probably "overttting"

No splitting of data in train, validation, test.

Already for our simple, p parameter in data points Lin. Reg. we have choices on the order/number of parameters. If our choices is inspired by the validation or test data, the ML model might not perform on new data.

Training data is used to set the weights:

$$\hat{w} = \arg \min_w \sum_{i \in \text{training}} w_i \text{loss} \left([y_i - f(x_i; w)] \right) + \text{reg}(w)$$

- Choices:

MLP
Multi layer percect...

- High level methods (in this course: lin. reg., k-means, neighbors)
- Integers, polynomials, cos/sin, choices within each method.
 - ↳ Not \mathbb{Z} , but in general yes/no questions.
 - ↳ especially neural network architecture
- loss functions
 - Robustness
 - Regularization
 - Symmetries
 - Sparseness
- Continuous parameters

↳ Hyperparameters, θ , are different than weights (continuous parameters), because they cannot be optimized. We have to make a choice.

Validation

$$\text{area } Q_\theta = \sum_{i \in \text{validation}} w_i \text{loss} \left(y_i - f(x_i; \hat{w}_\theta) \right)$$

↑
optimized on training

Multi-layer Perceptron

Perception: Linear action, non-linear response

: The small element that intelligent systems are made of.

$$y = f(x; w), \hat{w} = \arg \min_w \left[\sum_{i \in \text{train}} \text{loss}(y_i - f(x_i; w)) + \lambda \|w\| \right]$$

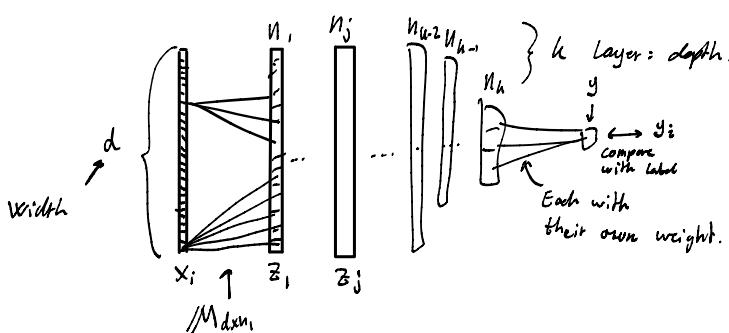
MLP ↗ Last time ↗ \hat{w} is the specific case for Lin. reg.

The simplest case, where x has a fixed shape:

$$x_i \in \mathbb{R}^d \quad (\text{Just a list of numbers})$$

$$X_{\text{train}} = [x_1, x_2, \dots, x_n]^T$$

\uparrow # of data samples



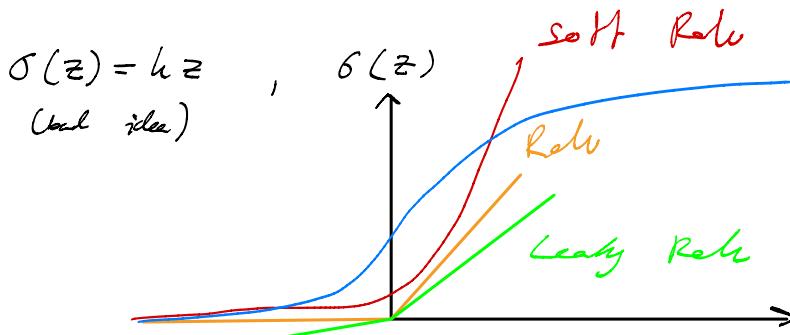
If every node in one layer is connected to every node in the next layer, the graph is said to be fully connected. If not, the network is said to be sparse.

Weights for each d entry connecting n_i nodes in the next layer. The weights depend on all other entries in a non-linear way?

$$z_j = \sigma(W_j \cdot z_{j-1} - b_j)$$

↑ $n_j \times n_{j-1}$ n_{j-1} n_j
 Non-lin. included here, applied to every entry of the
 vector in parenthesis. They're called Activation Functions

Examples of non-linearities:



Lecture 5

Optimization in Neural Networks

Call that in the Training we wanted:

$$\underset{w}{\operatorname{argmin}} \sum_{i \in \text{train}} \text{loss}(y_i - \hat{f}(x_i; w))$$

could be MLE

- Autodifferentiation: $\frac{df}{dw}|_{x_i}$ exists because the non-linearity is differentiable and its argument is simply linear.
- Back propagation

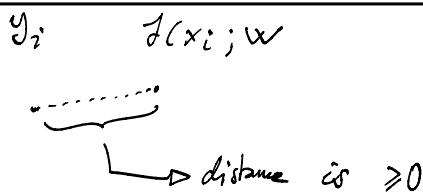
convex \leftarrow labels (contain cat?)
 total loss : $Q = \sum_{i \in \text{train}} \text{loss}(y_i - f(x_i, w))$
 ↗ data: (youtube videos)

$$\frac{d \text{loss}}{d w} = \text{loss}'(y_i - f(x_i, w)) \left[-\frac{\partial f}{\partial w} \Big|_{x_i} \right]$$

$$\Rightarrow \frac{d Q}{d w} = \sum_{i \in \text{train}} \frac{d \text{loss}_i}{d w}$$

A small (say 10%) subset of x_i can be used to become an estimator for $\frac{d Q}{d w}$.

convex function:



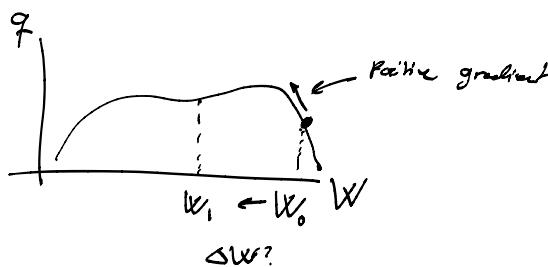
$$\hat{\frac{d Q}{d w}} = \frac{\text{Norm}}{N_{\text{sub}}} \sum_{i \in \text{sub}} \frac{d \text{loss}_i}{d w}$$

Stochastic gradient
Subset is randomly chosen.

↓
Batch

Leon Bottou: Stochastic gradient descent

Optimization:



Optimizers:

Newton's method: The next step, $\Delta w = -\frac{\frac{dq}{dw}}{\frac{d^2q}{dw^2}}$

Great for 1D w . For higher Dim. w :

- Conjugate gradient method (great)
- BFGS (Best, but complicated)

Why not Newton? Let's look at cost. For $\text{dim}(w) \approx 10^9$

$$\frac{\frac{dq}{dw}}{\frac{d^2q}{dw^2}}$$

$10^9 \times 10^9$ matrix. not diag.

For astronomical dim of w , we cannot use regular optimizers.

Stochastic Gradient Descent

You are at w_t at iteration t .

You know: $\frac{\hat{dq}}{dw} \Big|_{w_t}$

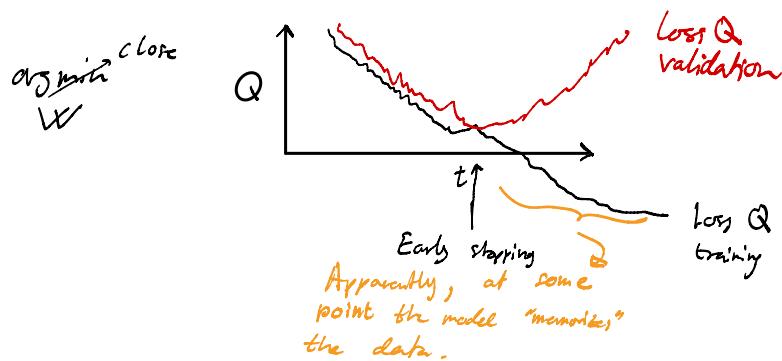
$$w_{t+1} = w_t + \alpha_t \frac{\hat{dq}}{dw} \Big|_{w_t}$$

Training rate.
 $\alpha_t = \frac{x_0}{(t+a)^s}$

Concurrent importance
 in ML. Might seem like a
 poor optimizer, but we cannot
 afford to take $\frac{dq}{dw}$ more than
 once.

For some choice
 of a , t
 you will be
 able to
 find the
 minimum.

- ① A better method: ADAM.
- ② The Landweber is terrible: There is only local minima.
Due to degeneracies, there will never be a global minimum
- ③ You don't want to be at the local minimum.



- ④ Performance on the validation is the same as test.

Lecture 6

Unsupervised ML

Data without labels, \mathbf{Y} .

Note:

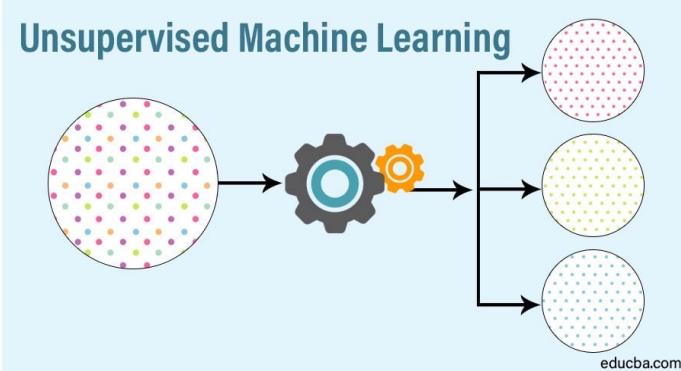
There is a small difference between self-supervised and unsupervised.

- ① Clustering \longleftrightarrow ② Low-Dim. structure.
- ③ De-noise data
- All these are related!
- Maybe the pixels pr data, but maybe it is only a subset that matters?
- ④ Predict or generate new data

Unsupervised Machine Learning

Clustering / Labelling :

- ⑤ Assess the consistency
of new data.



Methods of unsupervised learning:

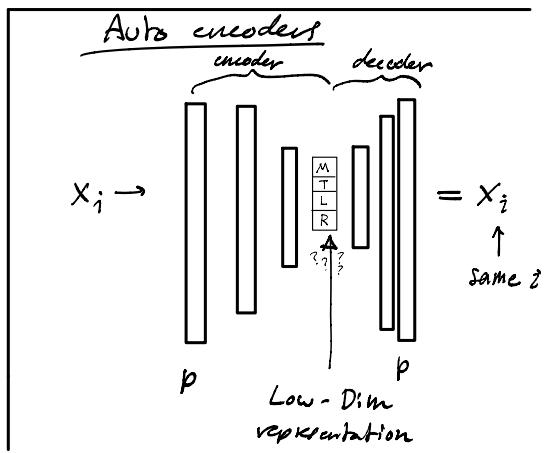
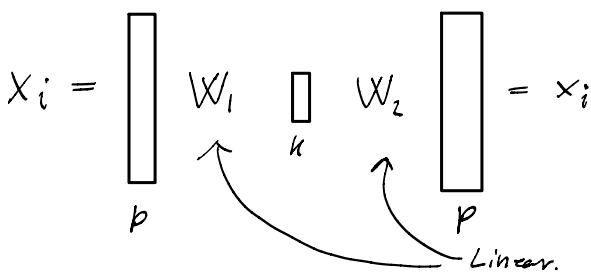
- ## (A) Principal Component Analysis

Features, x . Each x are a p -vector. All features in one:

$$X = [x_1, x_2, \dots, x_n]^T = A \cdot G + \text{residuals}$$

Matrix factorization
 $n \times p$ $n \times h$ $h \times p$
 $\underbrace{}$ rank k , $k < \min(n, p)$

$$Q = \underset{A, G}{\operatorname{arg\,min}} \|X - A \cdot G\|_2^2 \quad @ \text{rank } k.$$



METHODS FOR small Data:

Single Value Decomposition: SVD

"Finding eigenvalues" for rectangular matrices.

$$A = \begin{matrix} U & S \\ U^T & V \end{matrix}$$

$n \times p$ $n \times q$ $q \times q$ $q \times p$

$q = \min(n, p)$

Recall: Square X :

$$X = U^T X U$$

Recall pseudo inverse

$$(X^T X)^{-1} X^T = X^{-1} \quad n > p$$

$$X^T (X X^T)^{-1} = X^{-1} \quad n < p$$

Pseudo inverse for PCA:

$$X^T = V^T \frac{1}{S} V^T$$

2 for zero,
return zero.

Big Data: Low-Rank SVD

One liner in "PsyKitLearn".

Comments on PCA:

- ① Workhorse
- ② Outliers are bad: $X = A \cdot S + \underset{\uparrow}{\text{residual}} + \text{outliers}$
- ③ $p \rightarrow \infty$ in kernel PCA handles outliers.